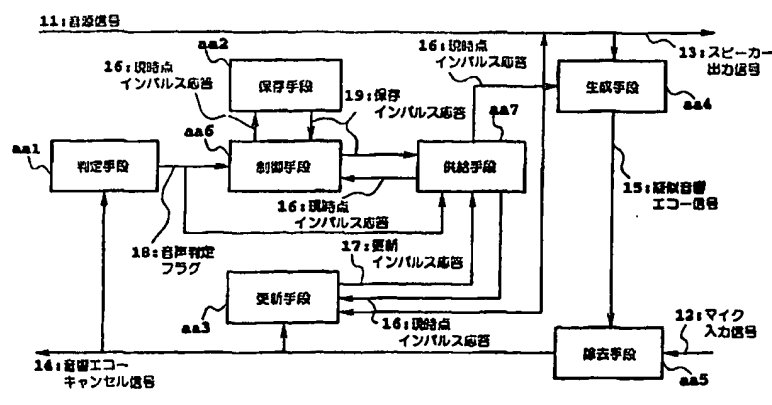




(51) 国際特許分類 H04R 3/02	A1	(11) 国際公開番号 WO98/39946  (43) 国際公開日 1998年9月11日(11.09.98)
(21) 国際出願番号 PCT/JP98/00915  (22) 国際出願日 1998年3月5日(05.03.98)  (30) 優先権データ 特願平9/51577 1997年3月6日(06.03.97) JP  (71) 出願人 (米国を除くすべての指定国について) 旭化成工業株式会社 (ASAHI KASEI KOGYO KABUSHIKI KAISHA)[JP/JP] 〒530-8205 大阪府大阪市北区堂島浜1丁目2-6 Osaka, (JP)  (72) 発明者; および (75) 発明者/出願人 (米国についてのみ) 庄境 誠(SHOZAKAI, Makoto)[JP/JP] 〒243-0216 神奈川県厚木市宮の里4-1-4-501 Kanagawa, (JP) 谷 智洋(TANI, Tomohiro)[JP/JP] 〒243-0021 神奈川県厚木市岡田4-5-3 厚木ユースハイム115号 Kanagawa, (JP)  (74) 代理人 弁理士 谷 義一(TANI, Yoshikazu) 〒107-0052 東京都港区赤坂5丁目1-31 第6セイコービル3階 Tokyo, (JP)	(81) 指定国 AL, AM, AT, AU, AZ, BA, BB, BG, BR, BY, CA, CH, CN, CU, CZ, DE, DK, EE, ES, FI, GB, GE, GH, GM, GW, HU, ID, IL, IS, JP, KE, KG, KR, KZ, LC, LK, LR, LS, LT, LU, LV, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, SL, TJ, TM, TR, TT, UA, UG, US, UZ, VN, YU, ZW, ARIPO特許 (GH, GM, KE, LS, MW, SD, SZ, UG, ZW), ユーラシア特許 (AM, AZ, BY, KG, KZ, MD, RU, TJ, TM), 欧州特許 (AT, BE, CH, DE, DK, ES, FI, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI特許 (BF, BJ, CF, CG, CI, CM, GA, GN, ML, MR, NE, SN, TD, TG).  添付公開書類 国際調査報告書	

(54)Title: DEVICE AND METHOD FOR PROCESSING SPEECH

(54)発明の名称 音声処理装置および方法



11 ... sound source signal  
12 ... microphone input signal  
13 ... loudspeaker output signal  
14 ... pseudo-acoustic echo signal  
15 ... pseudo-acoustic echo signal  
16 ... current impulse response  
17 ... updated impulse response  
18 ... speech judging flag  
19 ... stored impulse response  
aa1 ... judging means  
aa2 ... storage means  
aa3 ... updating means  
aa4 ... generating means  
aa5 ... removing means  
aa6 ... control means  
aa7 ... supplying means

## (57) Abstract

A speech processor which continuously uses the impulse response used by the preceding frame when a speech signal is included in the microphone input signal or the currently updated impulse response when no speech signal is included in the microphone input signal as the impulse response used for generating a pseudo-acoustic echo signal at the time of performing echo cancellation by using the pseudo-acoustic echo signal.

(57) 要約

擬似音響エコー信号を使用してエコーキャンセルを行う際に、擬似音響エコー信号の発生のために使用するインパルス応答として、マイク入力信号に音声が含まれる場合には前の時点のフレームで使用したインパルス応答を連続的に使用し、マイク入力信号に音声が含まれない場合には新規に更新されたインパルス応答を使用する音声処理装置。

PCTに基づいて公開される国際出願のパンフレット第一頁に掲載されたPCT加盟国を同定するために使用されるコード (参考情報)

AL	アルバニア	FI	フィンランド	LT	リトアニア	SN	セネガル
AM	アルメニア	FR	フランス	LUV	ルクセンブルグ	SZ	スワジランド
AT	オーストリア	GB	英国	LV	ラトヴィア	TD	トogo
AC	オーストラリア	BE	ベルギー	MC	モナコ	DJ	ジブチ
AZ	アゼルバイジャン	GE	グルジア	MD	モルドバ	TJ	タジキスタン
BA	ボスニア・ヘルツェゴビナ	GM	ガナ	MG	マダガスカル	TM	トルクメニスタン
BB	バルバドス	GN	ギニア	MK	マケドニア共和国	TR	トルコ
BE	ベルギー	GW	ギニア・ビサウ			TT	トリニダード・トバゴ
BG	ブルガリア	GR	ギリシャ	ML	マリ	UA	ウクライナ
BJ	ベナン	HU	ハンガリー	MN	モンゴル	UG	ウガンダ
BR	ブラジル	IE	アイルランド	MR	マリタニア	US	米国
BY	ベラルーシ	IL	イスラエル	MX	メキシコ	UZ	ウズベキスタン
CA	カナダ	IS	アイスランド	NE	ニジェール	VN	ベトナム
CC	中央アフリカ共和国	IT	イタリア	NL	オランダ	YU	ユーゴスラヴィア
CG	コンゴ共和国	JP	日本	NO	ノルウェー	ZW	ジンバブエ
CH	スイス			NZ	ニュージーランド		
CI	コートジボワール	KE	ケニア	PL	ポーランド		
CM	カメルーン	KG	キルギス	PT	ポルトガル		
CN	中国	KR	韓国	RO	ルーマニア		
CU	キューバ	KZ	カザフスタン	RU	ロシア		
CY	キプロス	LC	セント・ルシア	SE	スウェーデン		
CZ	チェコ	LL	リベリア	SG	シンガポール		
DE	ドイツ	LR	レソト	SI	スロベニア		
DK	デンマーク	LS	レソト	SK	スロバキア		
EE	エストニア			SL	シエラレオネ		
ES	スペイン						

## 明細書

## 音声処理装置および方法

## 技術分野

本発明は、リモートスピーカーからリモートマイクへの音響エコーをキャンセルすることにより通話品質の向上を計ったハンズフリー系の通話システム（テレビ会議システム、自動車電話）およびリモートスピーカーからリモートマイクへの回り込み音声をキャンセルすることにより音声認識性能の向上を計ったハンズフリー系の音声認識装置（カーオーディオ、カーナビゲーション、PCなど）に適用可能な音声処理装置および方法に関する。

## 背景技術

リモートスピーカーからリモートマイクへ回り込む音響信号は、しばしば音響エコーと呼ばれる。音響エコーを除去する技術（音響エコーキャンセラー）の用途は以下の2つである。

1) ハンズフリー系通話システム（テレビ会議システム、自動車電話）において、通話をしている相手に対し送出される音声の音質を向上させる。

リモートスピーカーから出力された相手側の音声は部屋の壁や窓ガラスに反射し、その部屋固有の音響特性の影響を受けて、リモートマイクに音響的に回り込む場合がある。この場合、相手にとっては自分の声がある時間遅れを伴って音響エコーとして戻ってくるため、聞きづらく話しづらいという不具合がある。従って、リモートマイクが集音した音声の中で、スピーカーから回り込んだ音響エコーをキャンセルして、残りの音声を通話をしている相手に送出することにより、上記の不具合を改善することが望

まれる。

2) ハンズフリー系の音声認識装置において音声認識率を向上させる。

例えば、自動車内においては、カーオーディオやカーナビゲーションのスピーカー出力音が上記音響エコーと同様にダッシュボードや窓ガラスに反射して音声認識用マイクに回り込み、それが非定常の加法性雑音として作用して、音声認識率が低下するという不具合がある。従って、音声認識用マイクが集音した音声の中で、スピーカーから回り込んだ音声をキャンセルして、残りの音声の認識を行うことにより、より高い音声認識性能を実現することが望まれる。

上記2つの用途においてはいずれも、リモートスピーカーからの出力音響の直接音および部屋の壁、ダッシュボードや窓ガラスなどで反射した反射音が常時リモートマイクに回り込む。ここでは、リモートスピーカーからリモートマイクへの直接音および反射音をまとめて音響エコーと呼ぶことにする。また、リモートスピーカー出力音から音響エコーが生成される経路を音響エコー生成経路と呼ぶことにする。

一般に、音響エコー生成経路の特性は、FIR (Finite Impulse Response) フィルターでモデル化できるが、部屋内の状況（人間の動作、人数などの要因）や自動車内の状況（人間の動作、人数、窓の開閉などの要因）により変化すると考えられる。音響エコー生成経路の特性の変化がほとんど起こらない場合には、あらかじめ最適なフィルター係数を求めておき、フィルター係数を固定して、音響エコーをキャンセルする方法で良いと思われる。しかしながら、音響エコー生成経路の特性の変化がいつ発生するかは、一般に予測が困難である。この場合、適応フィルターの利用により、最適なフィルター係数を動的に推定しながら、適応的に音響エコーのキャンセルを行う方法の採用が望ましい。

適応フィルターとは、観測信号が、音源既知の信号に対しあるインパルス応答を持つフィルターが畳み込まれて生成されたものであると仮定し、観測信号と（音源既知信号とフィルター係数の推定値との畳み込みにより計算される）疑似信号の差が0（ゼロ）になるように、フィルターの係数を動的に適応させるアプローチをいう。音響エコー生成経路を近似するFIRフィルターの係数とスピーカからの出力信号の畳み込みにより得られる信号を観測信号から引くことにより、音響エコーをキャンセルすることが可能である。適応フィルターのアルゴリズムとして、これまでにLMS（Least Mean Square error）[S. Haykin, "Adaptive Filter Theory," 2nd ed. Englewood Cliffs, NJ, Prentice-Hall, 1991]、NLMS（Normalized Least Mean Square error）[S. Haykin, "Adaptive Filter Theory," 2nd ed. Englewood Cliffs, NJ, Prentice-Hall, 1991]、APA（Affine Projection Algorithm）[尾関和彦, 南雲仁一, "アフィン部分空間への直交射影を用いた適応フィルター・アルゴリズムとその諸性質," 信学論, Vol.J67-A, No.2, pp.126-132, 1984.], RLS（Recursive Least Squares）[S. Haykin, "Adaptive Filter Theory," 2nd ed. Englewood Cliffs, NJ, Prentice-Hall, 1991]などが提案されている。特に、NLMSは、演算量が少なく、収束速度が既知の音源信号の大きさに依存しないため、広く採用されている。しかし、音声のような有色信号に対する、フィルター係数の収束速度が、APAやRLSに比べて遅いことが指摘されている。

時刻  $t$  におけるFIRフィルターの係数、FIRフィルターへの入力データ（既知の音源信号）をそれぞれ

$$h(t)=[h_1(t), h_2(t), \dots, h_M(t)]^T \quad (1)$$

$$x(t)=[x(t), x(t-1), \dots, x(t-M+1)]^T \quad (2)$$

で表現する。ここで、 $T$ は転置を示す。 $M$ はFIRフィルターの次数である。また、時刻 $t$ のマイク入力信号を $y(t)$ とすると、NLMSは、一般に以下の式で与えられる。

$$r(t)=h(t)^T x(t) \quad (3)$$

$$e(t)=y(t)-r(t) \quad (4)$$

$$h(t+1)=h(t)+\frac{\mu}{a+\|x(t)\|^2}x(t)e(t) \quad (5)$$

ここで、 $\|\cdot\|^2$ はベクトルのエネルギーを表す。 $\mu$ は、フィルター係数の更新速度を決定する定数（ステップゲインと呼ばれる）で、フィルター係数が収束するために、 $0<\mu<2$ を満たす必要がある。 $a$ は、 $\|x(t)\|^2$ が微小値の場合に(5)式の右辺第2項が発散するのを防止するための正の定数である。図1に上述の式を回路で表したNLMSのブロック図を示す。ここで、 $r(t)$ を疑似音響エコー信号、 $e(t)$ を音響エコーキャンセル信号と呼ぶことにする。図2に、NLMSをはじめとする適応フィルターを使用した音響エコーキャンセラー（AEC）を室内に設置した例を示す。説明の便宜上、スピーカー2から出力される信号のAEC1への入力を遠端入力、マイク3の入力を近端入力、スピーカー2の出力を近端出力、音響エコーキャンセル後のAEC1の出力信号を遠端出力と呼ぶ。また、遠端入力と近端出力は全く等価であると仮定し、遠端入力から近端出力が生成される系の特性（スピーカー特性など）は、音響エコー生成経路の特性に含めるものとする。

このような音響エコーキャンセラーに関しては、特に、以下の課題について精力的に研究が進められてきた。

### 1) ステップゲインの制御

ステップゲインは可能な限り大きくして収束速度を上げることが必要であるが、大きくしすぎるとハウリングの原因となるため、使用環境に適した設定が必要である。代表的なステップゲインの制御方法として、E S (Exponential Step)法[S. Makino, Y. Kaneda and N. Koizumi, "Exponentially Weighted Stepsize NLMS Adaptive Filter Based on the Statistics of a Room Impulse Response" IEEE Trans. SAP, Vol.1, No.1, pp.101-108, 1993.]が提案されている。室内における、適応フィルター係数更新時の変化量が指数減衰特性を有することから、ステップゲインを指数的に(変化量の大きいインパルス応答前半では大きく、後半は小さく)設定する。残留エコーレベルが減少するのに要する時間が通常のNLMSの半分程度で済むことが示されている。

### 2) ダブルトーク検出

一般に、ダブルトーク(遠端話者と近端話者の双方が同時に発声した状態)において、AEC(NLMS)1により適応フィルター係数の更新を継続すると、フィルターの係数が大きく乱れ、その結果としてエコー消去量が減少してハウリングが起りやすくなる。従って、ダブルトークを如何に速やかに検出して、AEC1の適応フィルターの更新動作を制御(停止・再開)するかは重要なポイントである。ダブルトークの検出には、残留エコーのパワーを用いる方法が有効とされる[藤井健作, 大賀寿郎, "エコー経路変動検出を併用するダブルトーク検出法," 信学論, Vol.J78-A, No.3, pp.314-322, 1995.]. 音響エコーに埋もれる小さな、近端話者音声を検出することが可能であるからである。

### 3) 音響エコー経路変動検出

近端話者が移動した時などに伴う音響エコー経路の変動が発生した場合、

残留エコーが増大してダブルトークと判断され、適応フィルターの係数更新が停止されるという不具合が生じる。そこで、ダブルトークと音響エコー経路変動を区別し、音響エコー経路変動の場合には適応フィルターの係数更新を継続することが必要になる。その方法として、近端入力と疑似エコーの相関を利用する方法[藤井健作, 大賀寿郎, "エコー経路変動検出を併用するダブルトーク検出法," 信学論, Vol.J78-A, No.3, pp.314-322, 1995.]などが提案されている。

カーオーディオやカーナビの音響・音声信号がスピーカーから出力されている自動車環境内でのロバストな音声認識機能を高める目的で、車室内にNLMSによるAEC1を設置して音源既知の加法性雑音の除去する例を図3に示す。図3に示す図2と同一個所の符号は図2に示す符号と同一である。

スピーカー2からガイダンス音声出力されている最中に発声された音声認識できる機能、いわゆるBarge-In(Talk-Through)機能の実現に向けて、AEC1を利用する手法が試みられている。ここで、スピーカー2の出力に起因する音声の誤認識の中で、適応フィルターの効果により正認識となる回復率をRRE(Recovery Rate of Error)と呼ぶことにする。

例えば、会議室に設置された音声対話システムにおいて、スピーカー2からマイク3へのガイダンス音声の回り込みをAEC1により抑圧することにより、70～80%のRREが得られることが示されている[高橋敏, 嵯峨山茂樹, "NOVO合成法を用いたBarge-in音声の認識," 音講論集, 2-5-1, pp.59-60, 1996-3.]。

しかしながら、音源未知の加法性雑音が存在し、その雑音レベルが常時変動する車室内における、音響エコーキャンセラーに関する研究成果はあまり報告されていない。自動車電話のハンズフリー装置においては、音声



スイッチ（近端入力と遠端入力のエネルギー比較による交互通話方式）と音響エコーキャンセラーとの併用によるものもあるが、語頭、語尾の切断が多く通話品質が不十分である点が指摘されている。

一般に、近端入力に近端出力から生成される音以外の音が混入する（以下、近端入力が存在するという）状況で係数の適応化を継続した場合、フィルター係数の推定精度が劣化し、音響エコーのキャンセル性能が悪化する。そこで、遠端入力が存在し、かつ近端入力が存在する状態（ダブルトーク状態と呼ぶ）では、(5)式によるフィルター係数の更新を停止させることが一般に行われる。遠端入力が存在するかどうかの判断は、遠端入力のエネルギーと予め定められたしきい値との単純な比較で可能である。

一方、近端入力が存在するかどうかの判断を同様に行うと、音響エコーの影響で近端入力が存在すると判断するケースが多くなり、(5)式によるフィルター係数の更新を頻繁に停止して、結果的にフィルター係数の推定精度が劣化するという不具合が生じる。そこで、近端入力信号  $y(t)$  ではなく、音響エコーキャンセル信号  $e(t)$  のエネルギーを用いて、近端入力が存在するかどうかを判断するという方法が考えられる。近端出力から生成された音以外で近端入力に混入する音としては、大きく分けて走行雑音などの音源未知の加法性雑音と人間の音声の2つが考えられるが、いずれも適応フィルターで除去されずに遠端出力に残存する。

一般に、走行中の自動車環境では、音源未知の加法性雑音のエネルギーレベルが、60～80 dBAの間で大きく変動するため[金指久則，則松武志，新居康彦，"車載用単語音声認識装置," 音講論集，1-Q-32, pp.159-160, 1995-3.][鈴木邦一，中村一雄，宇尾野豊，浅田博重，"車載騒音環境下における連続音声認識," 音講論集，2-Q-4, pp.155-156, 1993-10.]、近端入力の存在を判断するための最適なしきい値を一意に決めることは難しいという問題

がある。

また、音源未知の加法性雑音の影響により、近端入力と疑似音響エコー信号の相関が低下するため、先述の音響エコー経路変動検出法[藤井健作, 大賀寿郎, "エコー経路変動検出を併用するダブルトーク検出法," 信学論, Vol.J78-A, No.3, pp.314-322, 1995.]の適用も困難な場合があると予想される。音源未知の加法性雑音と人間の音声を正確に識別する能力を持つ、音声検出アルゴリズムがあれば、有力な解決法になると思われる。

まず走行雑音のみが存在する場合での、NLMSによる音響エコーのキャンセル性能を評価する。図4A、図4B、図4C、図4D、図4Eに、それぞれ遠端入力信号（ポップス音楽）のスペクトログラム、アイドリング時での近端入力信号のスペクトログラム、同じく音響エコーキャンセル信号のスペクトログラム、時速100km走行時の近端入力信号のスペクトログラム、同じく音響エコーキャンセル信号のスペクトログラムを示す。

カーオーディオの音量は、アイドリング時と時速100km走行時で、男性1名が快適と感じるレベルにセットした。従って、時速100km走行時の方が、スピーカー出力レベルは大きく、音響エコーレベルも大きい。近端入力信号は、2000ccの自動車の運転席サンバイザーに単一指向性マイクを設置して収録した。フィルター係数の初期値は全て0.0とし、時刻0秒から継続的に(3)-(5)式によりフィルター係数を更新しながら音響エコーキャンセル信号を求めた。サンプリング周波数は8kHzであり、音響エコーの最大遅延は32msまで考慮した。従って、FIRフィルターのタップ数は256である。

また、適応フィルターの性能を評価する尺度として、ERLE(Echo Return Loss Enhancement)がよく用いられる。ERLEは近端入力信号の減衰量を表し、次式で定義される[北脇信彦編著, "音のコミュニケーション工

学—マルチメディア時代の音声・音響技術—," コロナ社, 1996.]。

$$ERLE = 10 \cdot \log_{10} \frac{E[y(t)^2]}{E[e(t)^2]} \quad (6)$$

$E[\cdot]$ は推定値を表し、次式により求める。

$$E[z(t)^2] = (1 - \lambda) \cdot E[z(t-1)^2] + \lambda \cdot z(t)^2 \quad (7)$$

但し、 $\lambda = 1/256$ である。ERLEの単位は、dBである。アイドリング時のERLEの最大値、平均値はそれぞれ18.80 dB、10.13 dBである。また、時速100 km走行時のERLEの最大値、平均値はそれぞれ9.33 dB、5.89 dBである。近端入力音源未知の加法性雑音のレベルが大きいほど、(6)式で与えられるERLEは低い値になることに注意する必要がある。

図4C、図4Eからアイドリング時、時速100 km時いずれの場合も音響エコーをほぼキャンセルできていることが分かる。近端入力に人間の音声が含まれない場合は、フィルター係数を継続的に更新することにより音響エコーの大部分はキャンセル可能であると思われる。すなわち、音源未知の加法性雑音の中で定常的かつ音声と無相関である走行雑音は、フィルター係数の推定に与える影響が小さいと考えられる。

次に、近端入力に人間の音声が含まれる場合について調べる。2000 ccの自動車からカーオーディオからポップス音楽を再生しながら市街地を時速60 kmで走行し、加法性雑音データを収録した。この時、音楽のボリュームは女性1名が快適と感じるレベルにセットした。次に、停止中（エンジンオフ）の同一の自動車内で同一女性1名が発声した音声データ（「明るい」）を同一の録音レベルで収録した。

そして、加法性雑音データと音声データとを計算機上で加算した信号のスペクトログラムを図 7 A に示す。図 7 B にフィルター係数の初期値を 0.0 とし、時刻 0 秒から連続的にフィルター係数を更新した場合の音響エコーキャンセル信号のスペクトログラムを示す。また、図 7 C にフィルター係数の 10 番目の係数の値の変化を示す。この時の、E R L E の最大値、平均値はそれぞれ 8.48 dB、4.18 dB である。

特に、時刻 0.5 秒あたりから 0.15 秒の間、フィルター係数値が激しく振動し、不安定になっている様子が分かる。また、時刻 1.0 秒以降の音響エコー（図 7 B の楕円で囲まれた部分）をキャンセルできていない。近端入力に音声が存在する間はフィルター係数の更新を停止し、近端入力に音声が存在しない間は、定常的な加法性雑音の存在の如何に関わらずフィルター係数の更新を継続する必要がある。そのためには、音源未知の加法性雑音が混入する近端入力に音声が含まれているかどうかを正確に判定する音声検出アルゴリズムが必要となる。

音声認識システムにおいては、正確に音声区間を検出すること（音声検出）が極めて重要である。背景雑音がほとんどない環境では、正確な音声検出はそれほど難しくはない。しかしながら、走行中の車室内のように背景雑音の存在が無視できない環境においては、音声の検出はかなり困難である。特に、音声の最初に位置する弱い摩擦音、弱い鼻音や音声の最初または最後に位置する無声化した母音などは背景雑音に埋もれてしまうケースが多く、検出は難しい。呼吸音、舌打ち音などは本来非音声として検出されるべきであるが、しばしば音声として検出され、誤認識につながることが多い。

通常、あるしきい値以上の短時間パワーが連続して一定フレーム以上継続するかどうかにより音声の開始点を検出し、あるしきい値以下の短時間

パワーが連続して一定フレーム以上継続するかどうかにより音声の終了点を検出する方法が一般的である。また、2つのレベルのしきい値を用いて、より正確に音声を検出しようとする試みや、音声信号の零交差回数を用いることもある[古井貞熙, "デジタル音声処理," デジタルテクノロジーシリーズ, 東海大学出版会, 1985.]。音源未知の加法性雑音の存在が無視できる環境においては、短時間パワーや零交差回数などの時間情報のみを用いる音声検出法でも問題は生じない。

しかし、音源未知の加法性雑音の存在が無視できない環境においては、従来の音声検出法を用いた場合の音響エコーキャンセラーには、以下の不具合がある。まず、第一にマイク入力に音声が存在しないにも関わらず、音源未知の加法性雑音を音声であると判断し、フィルター係数の更新が行われず、音響特性の変化に追従できなくなり、音響エコーのキャンセル性能が低下する。第2に、マイク入力に音声が存在するにも関わらず、音声がないと判断し、フィルター係数の更新が行われて、所望の値から乖離し、音響エコーのキャンセル性能が低下する。従って、時間情報ばかりではなく、スペクトルなどの周波数情報も併用する方法が望ましい。

特開平9-213946号(NTT)においては、入力音声信号(エコーキャンセル前の信号)と音源既知の加法性雑音の音源情報の時間情報および周波数情報を用いて入力音声信号に音声が含まれているかどうかを判定するダブルトーク検出回路を用いた音響エコーキャンセラーが説明されている。しかしながら入力音声信号に入り込むエコーとして音源信号の影響を受けたものだけを前提にしており、周囲の雑音がある場合に、ダブルトーク検出精度が悪いという不具合がある。また、適応フィルターにより推定したインパルス応答(FIRフィルターの係数値)を保持するバッファを有していない。

特開平５－１０２８８７号（東芝）では、エコーキャンセル後の信号の大きさにダブルトークかどうか判定するダブルトーク検出回路を用いているが、時間情報および周波数情報を併用する判定ではないため、周囲の雑音が存在する環境における判定精度が十分ではないという不具合がある。

特開平７－３０３０６６号（NTT DOCOMO）では、インパルス応答レジスタで判定手段の遅れを補償する構成を取っているが、エコーキャンセル後の信号の時間情報および周波数情報を用いて入力音声信号に音声が含まれているかどうかをフレーム毎に判定する手段を具えていないため、ダブルトーク検出性能に限界がある。

WO 96/42142号（NOKIA）では、エコーキャンセル後の信号の時間情報および周波数情報を用いて入力音声信号に音声が含まれているかどうかをフレーム毎に判定する手段を具えているが、自動車電話の基地局の送出信号のゲインを小さくすることにより音響エコーが直接送出されるのを押さえる構造を持つ音響エコーサプレッサに関する発明であり、音響エコーキャンセラーに関する発明ではない。

## 発明の開示

本発明の目的は、音響などの雑音が混在しやすい環境下で音声信号からの雑音除去性能を改善することができる音声処理装置および方法を提供することにある。

音響エコーの伝達経路を模擬する現時点のインパルス応答および音源信号に基づき疑似音響エコー信号を生成する生成手段と、

現時点のインパルス応答を保持し、前記生成手段に供給する供給手段と、

マイク入力信号から該疑似音響エコー信号を減算することにより音響

エコー成分を除去し、音響エコーキャンセル信号を生成する除去手段と、

前記音源信号と前記音響エコーキャンセル信号と前記供給手段が保持している現時点のインパルス応答を用いて継続的にインパルス応答を更新し、更新されたインパルス応答を前記供給手段に供給する更新手段と、

前記音響エコーキャンセル信号の時間情報および周波数情報を利用して、マイク入力信号に音声が含まれているか否かをフレーム毎に判定する判定手段と、

1つ以上のインパルス応答を保存する保存手段と、

前記判定手段の判定結果が否定判定のフレームでは前記供給手段が保持している現時点のインパルス応答を前記保存手段に保存し、肯定判定のフレームでは前記保存手段に保存されているインパルス応答の1つを取り出して、前記供給手段に供給する制御手段と

を具えたことを特徴とする。

本発明では前記音響エコーがキャンセルされた後の信号を音声認識に用いてもよい。

本発明ではさらに前記音響エコーがキャンセルされた後の信号から、フーリエ変換により各フレーム毎にスペクトルを求める手段と、当該得られたスペクトルに基づき各フレーム毎に連続的にスペクトル平均を求める手段と、当該得られたスペクトル平均を前記音響エコーがキャンセルされた後の信号から各フレーム毎に計算されたスペクトルから連続的に減算することにより、音源未知の加法性雑音を除去する手段とを具えてもよい。

本発明ではさらに前記音響エコーがキャンセルされた後の信号から、フーリエ変換により各フレーム毎にスペクトルを求める手段と、当該得られたスペクトルから各フレーム毎に連続的にスペクトル平均を求める手段と、当該得られたスペクトル平均を前記音響エコーがキャンセルされた後の信

号から各フレーム毎に計算されたスペクトルから連続的に減算することにより、音源未知の加法性雑音を除去する手段と、当該加法性雑音が除去されたスペクトルからケプストラムを求める手段と、当該得られたケプストラムの音声フレームのケプストラム平均および非音声フレームのケプストラム平均を話者毎に別々に求め手段と、話者毎に音声フレームのケプストラムからはその話者の音声フレームのケプストラム平均を減算し、非音声フレームのケプストラムからはその話者の非音声フレームのケプストラム平均を減算して、マイク特性や口からマイクまでの空間伝達特性に依存する乗法性歪みを補正する手段とを具えてもよい。

本発明ではさらに、前記音響エコーがキャンセルされた後の信号から、フーリエ変換により各フレーム毎にスペクトルを求める手段と、当該得られたスペクトルからケプストラムを求める手段と、当該得られたケプストラムの音声フレームのケプストラム平均および非音声フレームのケプストラム平均を話者毎に別々に求める手段と、話者毎に音声フレームのケプストラムからはその話者の音声フレームのケプストラム平均を減算し、非音声フレームのケプストラムからはその話者の非音声フレームのケプストラム平均を減算することにより、マイク特性や口からマイクまでの空間伝達特性に依存する乗法性歪みを補正する手段とを具えてもよい。

本発明では、フーリエ変換により各フレーム毎にスペクトルを求める手段と、当該得られたスペクトルからケプストラムを求める手段と、当該得られたケプストラムの音声フレームのケプストラム平均および非音声フレームのケプストラム平均を話者毎に別々に求める手段と、話者毎に音声フレームのケプストラムからはその話者の音声フレームのケプストラム平均を減算し、非音声フレームのケプストラムからはその話者の非音声フレームのケプストラム平均を減算することにより、マイク特性や口からマイ



クまでの空間伝達特性に依存する乗法性歪みを補正する手段とを具えてもよい。

本発明では、擬似音響エコー信号を使用してエコーキャンセルを行う際に、擬似音響エコー信号の発生のために使用するインパルス応答として、マイク入力信号が音声の場合には前の時点のフレームで使用したインパルス応答を連続的に使用し、マイク入力信号が音声ではない場合には新規に更新されたインパルス応答を使用することで音響エコーキャンセリングの性能を改善する。

さらに本発明は音響エコーをキャンセルした後の信号からフレーム毎のスペクトルおよびスペクトル平均を求め、得られたスペクトルおよびスペクトル平均を使用して加法性雑音を除去する。

#### 図面の簡単な説明

図 1 は N L M S (Normalized Least Mean Square error) の機能構成を示すブロック図である。

図 2 は音響エコーキャンセラーの設置例を示す図である。

図 3 は車室内における音源既知の加法性雑音を除去する例を示す図である。

図 4 A - 図 4 E はそれぞれ N L M S (Normalized Least Mean Square error) の性能 (横軸: s e c . ) を示す図である。

図 5 は V A D (Voice Activity Detection) の処理内容を示すブロック図である。

図 6 は V A D の動作タイミングを示す図である。

図 7 A - 図 7 G はそれぞれ N L M S - V A D (Normalized Least Mean Square error with frame-wise Voice Activity Detection) の効果 (横軸: s e

c.)を示す図である。

図8はフィルター係数バッファの動作を説明するための図である。

図9はNLMS-VADの構成を示すブロック図である。

図10Aおよび図10BはそれぞれNLMS-VAD/CSS法によるスペクトログラムを示す図である(横軸:sec.)。

図11は時不変フィルタを示す図である。

図12はNLMS-VAD/CSS/E-CMN法の処理内容を示すブロック図である(横軸:sec.)。

図13はNLMS-VAD/CSS/E-CMNの評価を示す図である。

図14は本発明第1実施形態の音声処理装置の構成を示すブロック図である。

図15は本発明第2実施形態のシステムの構成を示すブロック図である。

図16は本発明第3実施形態のシステムの構成を示すブロック図である。

図17は本発明第4実施形態のシステムの構成を示すブロック図である。

図18は本発明第5実施形態のシステムの構成を示すブロック図である。

図19は本発明第6実施形態のシステムの構成を示すブロック図である。

#### 発明を実施するための最良の形態

短時間パワーやピッチなどの時間情報とスペクトルなどの周波数情報を利用した音声検出アルゴリズムの1つとして、欧州の携帯電話・自動車電話システムであるGSMで標準規格化されている音声検出VAD(Voice Activity Detection)[Recommendation GSM 06.32.]がある。このVADは音声CODEC(圧縮・伸張)などのデジタル信号処理の動作を細かく制御し、低消費電力化を計って電池寿命を延ばす目的で用いられている。図5にこのVADの簡単な構成を示す。まず、音声信号からフレーム毎に自己

相関関数（時間情報）が求められる。この自己相関関数から線形予測分析 L P C（Linear Predictive Coding）により、線形予測係数（時間情報）が求められる。線形予測係数から構成できる逆 L P C フィルターと自己相関関数から音声信号の短時間パワー（時間情報）を求めることができる。この短時間パワーとしきい値を比較し、V A D 判定を行う。

短時間パワーがしきい値よりも大きい場合は、値 1 の局所的な V A D フラグが出力される。そうでない場合は、値 0（ゼロ）の局所的な V A D フラグが出力される。そして、V A D 後処理において過去の複数のフレームの局所的な V A D フラグの値の履歴を用いて最終的な V A D フラグの値が決定される。

一方、V A D 判定において短時間パワーとの比較に用いられるしきい値は、以下のように適応化される。平滑化された自己相関関数と自己相関予測係数により表されるスペクトル（周波数情報）変化が連続するフレーム間で十分小さい場合は、スペクトルの定常性が高いと判断される。スペクトルの定常性が高いと判断される音声信号としては背景雑音または母音と考えられる。

背景雑音のフレームにおいてはしきい値適応を行い、母音のフレームではしきい値適応を行うべきではない。背景雑音と母音を区別するためにピッチ情報を利用する。音声信号から計算される自己相関関数からピッチラグ（ピッチ周期）（時間情報）が計算される。連続するフレーム間でピッチラグの変化が小さい場合は、そのフレームは母音であると判断され、値 1 のピッチフラグが出力される。そうでない場合は、値 0（ゼロ）のピッチフラグが出力される。

上記の短時間パワー、逆 L P C フィルターから求められる残差信号自己相関予測係数、ピッチフラグ、定常性の情報を利用してスペクトルの定常

性が高く、ピッチ性が低いフレームにおいてしきい値の適応化が行われる。このVADはエンジン音やロードノイズなどの比較的定常的な背景雑音にたいしてはそのレベルに関わらず、正確な音声検出性能を発揮する。

自動車内において、安全性の見地からマイクがサンバイザーなど口元から離れた場所に設置される場合、信号雑音比（SNR）は10 dB以下にまで悪化する。その場合には、上記のVADアルゴリズムの音声検出性能は著しく劣化することが分かった。

そこで、SNR 10 dB程度でも正確に音声検出ができるようにしきい値の適応化などの部分を改良した。現在、VADに用いている音声の窓長は32 ms、フレームシフトは10 msである。以後、VADにより音声の存在が検出された場合、VADがONであると言う。逆に、検出されなかった場合、VADがOFFであると言う。このVADは、1フレームに1回近端入力に音声が含まれているかどうかの判断を行うため、音声の検出タイミングが実際の音声の立ち上がりから遅延することが起こりえる。

図6に、フレーム、VADの動作タイミングおよびVADが使用する窓長の関係を示す。実際の音声の開始がフレームnの中心である場合、VADによりその音声の開始を検出できるのは、フレームn+1以降である可能性が高い。仮に、フレームn+2のVADで検出できた場合、実際の音声の開始からの検出遅れは25 msにもなり、その間、エコーパスの推定値が不安定になることが考えられる。

図7DにVADによる音声検出の様子を示す。レベル1が音声を検出されたことを示す。矢印を用いて示すように、2フレーム程度の音声検出遅れが認められる。不安定になったフィルター係数値をより精度の高い値に回復することができれば、音響エコーキャンセル性能の低下を避けることが可能だと考えられる。

そこで、 $m$ 個分のフィルター係数を格納できるバッファ（フィルター係数バッファと呼ぶ）を用意する。VADがOFFのフレームでは、 $n$  ( $m-1 \geq n \geq 1$ ) 番目の格納位置に格納されたフィルター係数を順次  $n+1$  番目の格納位置に移すと同時に、現時点での適応フィルターの係数をフィルター係数バッファの第1番目の格納位置に格納する。この時、結果として、 $m$ 番目の格納位置に格納されていたフィルター係数は捨てられる。一方、VADがONのフレームでは、フィルター係数バッファの  $m$  番目の格納位置に格納されたフィルター係数を取り出し、その値で劣化したフィルター係数をリセットすれば良い。

図8にフィルター係数バッファの動作の様子を示す。 $m$ を0～4とした場合のERLEの最大値、平均値を表1に示す。

表1 フィルターバッファのサイズとERLE (Echo Return Loss Enhancement) の関係

buffer size $m$	max ERLE (dB)	average ERLE (dB)
0	8.80	4.18
1	9.06	4.25
2	9.15	4.35
3	9.14	4.36
4	9.14	4.36

$m = 0$  は係数値の保存およびリセットを行わない場合を示す。 $m \geq 2$  では、E R L E にほとんど違いが見られないため、 $m = 2$  を選択する。これは、V A D の検出遅れ（2 フレーム程度）と対応している。

上記の特徴を持ったアルゴリズムを N L M S - V A D (NLMS with frame-wise VAD) と呼び、全体のブロック図を図 9 に示す。ここで、[s]、[f] はそれぞれサンプルワイズ、フレームワイズの信号の流れおよび処理の動作を示す。V A D がいったん O N になると、次に V A D が O F F になるフレームまでフィルター係数の更新は停止される。図 7 E に、フィルター係数の初期値を全て 0. 0 とし、V A D を動作させ、フィルター係数値の格納およびリセットを行いながら、時刻 0 秒からフィルター係数を更新した場合の音響エコーキャンセル信号のスペクトログラムを示す。図 7 F にその時のフィルター係数の 10 番目の係数の値の変化を示す。フィルター係数の更新が停止されたフレームの直前で、フィルター係数値が不安定になっているが、上記フィルター係数の格納およびリセットにより、フィルター係数が回復されている様子が示されている。これにより、時刻 1. 0 秒以降の音響エコー（図 7 B の楕円で囲まれた部分）もキャンセルされている。

但し、図 7 E で時刻 0. 1 秒前後の音響エコー（図 7 E の楕円で囲まれた部分）がキャンセルされていないことが分かる。発声毎に推定されたフィルター係数および V A D に用いられるパラメータを保存しておき、次の発声時にそれらを初期値として用いれば、フィルター係数の推定速度は速まると考えられる。図 7 G にその例を示す。時刻 0. 0 秒直後の音響エコーは若干残存しているが、それ以後の音響エコー（図 7 E の楕円で囲まれた部分）はほぼキャンセルされていることが分かる。この時の、E R L E の最大値、平均値はそれぞれ 9. 29 dB、4. 50 dB である。また、本

願出願人はNLMS-VAD法に関連して、時間情報および周波数情報に基づく音声検出を用いた音響エコーキャンセラーの日本国出願を既に済ませている（特願平09-051577号、1997年3月6日出願）。なお、本願発明では時間情報および周波数情報に基づく音声検出をフレーム毎に行う点が、上記先願発明との相違点である。

次に、音源既知の加法性雑音および音源未知の加法性雑音が存在する環境におけるロバストな音声認識方法として、NLMS-VAD法とCSS (Continuous Spectral Subtraction)法を組み合わせる方法について説明する。時刻  $t$  における周波数  $\omega$  での観測スペクトル、音声スペクトル  $S(\omega; t)$  の推定値、加法性雑音の推定値をそれぞれ  $O(\omega; t)$ 、 $\hat{S}(\omega; t)$ 、 $\hat{N}(\omega; t)$  と表す とすると、CSS法は以下のように与えられる。

$$\hat{N}(\omega; t) = \gamma \cdot \hat{N}(\omega; t-1) + (1-\gamma) \cdot O(\omega; t) \quad (8)$$

$$\hat{S}(\omega; t) = \begin{cases} O(\omega; t) - \alpha \cdot \hat{N}(\omega; t) & \text{if } O(\omega; t) - \alpha \cdot \hat{N}(\omega; t) > \beta \cdot O(\omega; t) \\ \beta \cdot O(\omega; t) & \text{otherwise} \end{cases}$$

(9)

ここで、 $\alpha$ はover-estimation factor、 $\beta$ はflooring factor、 $\gamma$ はsmoothing factorであり、以下では予備実験の結果から、それぞれ2.4、0.1、0.974と設定した。CSSは、音声フレームと非音声フレームを区別せず、連続的にスペクトルの移動平均を求め、これを雑音スペクトルの推定値とみなして、入力スペクトルから減算する方法である。雑音スペクトルの推定値に音声スペクトルの影響が含まれるため、エネルギーの弱い音声スペクトルがマスクされてしまい、歪みが生じるという問題点があるが、過去のある一定時間長の区間に対して、相対的に大きなエネルギーを持つ周波

数成分を残し、エネルギーの微弱な周波数成分を雑音、音声を問わず、マスクするという働きを持つ。このため、クリーンな音声にCSSを施した後に得られる特徴パラメータと加法性雑音が重畳した音声にCSSを施した後に得られる特徴パラメータの間の変動が、通常のスペクトル減算法や最小平均二乗誤差推定法に比べて小さい。この特長は、低いSNRでの音声認識にとって有効である。図10Aに、停止中（アイドリング）の自動車内で女性が発声した音声（「明るい」、図7Aに示した音声を計算機上で加算して作成した際に用いた音声と同一）にCSSを施した後のスペクトログラムを、図7Bに同一音声に時速60km走行時の音源未知の加法性雑音と音響エコーが重畳した雑音データを計算機上で加算した後（図7A）、NLMS-VAD法で音響エコーをキャンセルし（図7G）、CSS法を施して得られるスペクトログラムを示す。図7Gと図10Bを比較すると、時刻0.9秒近辺の周波数1kHzの音響エコーの残存成分（図7Gの楕円で囲まれた部分）がCSS法により除去されていることが分かる。

CSS法は、定常的な加法性雑音だけでなく、NLMS-VAD法でキャンセルできなかった残存音響エコーを抑圧する効果も持っている。音響エコーキャンセル信号 $e(t)$ にFFTを施して得られたスペクトルに対してCSS法を施した後のスペクトルを逆FFTにより時間領域に戻して得られる波形信号を、(6)式の $e(t)$ の代わりに用いた場合のERLEの平均値は13.60dBであった。これに対し、NLMS-VAD法による音響エコーキャンセルを行わず、CSS法による加法性雑音のキャンセルのみを行って同様に求めたERLEの平均値は9.87dBであった。CSSのみでは、約3.7dB相当の音源既知の加法性雑音がキャンセルできなかったと見ることができる。



図10Aと図10Bを比較すると2つのスペクトログラムがきわめて類似していることが分かる。NLMS-VAD法とCSS法の組み合わせにより、音源既知の加法性雑音と音源未知の加法性雑音に対して、ロバストな特徴パラメータを抽出できることが示唆されている。

次に、音声スペクトルに対する乗法性歪みの補正方法について述べる。ある個人の発声器官で生成される、時刻  $t$  における周波数  $\omega$  での短時間スペクトル  $S(\omega; t)$  の音声フレームにおける長時間平均を話者の個人性  $H_{person}(\omega)$  と呼ぶこととし、

$$H_{person}(\omega) = \frac{1}{T} \cdot \sum_{t=1}^T S(\omega; t) \quad (10)$$

と定義する。ここで、 $T$  は十分大きな自然数である。 $H_{person}(\omega)$  は、声帯音源特性および声道長に依存する話者固有の周波数特性を表しているとみなすことができる。また、短時間スペクトルを話者の個人性で除したもの

$$S^*(\omega; t) = S(\omega; t) / H_{person}(\omega) \quad (11)$$

を正規化音声スペクトルと定義する。この時、図11に示すように、音声スペクトルは、正規化音声スペクトル  $S^*(\omega; t)$  が時不変フィルター  $H_{person}(\omega)$  を通過することにより生成される、あるいは、正規化音声スペクトル  $S^*(\omega; t)$  に乗法性歪み  $H_{person}(\omega)$  が重畳して生成されると解釈することができる。

$$S(\omega; t) = H_{person}(\omega) \cdot S^*(\omega; t) \quad (12)$$

車室内のような実環境においては、正規化音声スペクトルに対する乗法性歪みとして、上記の話者の個人性に加えて、以下の3種類が考えられる [A. Acero, "Acoustical and Environmental Robustness in Automatic Speech Recognition," Kluwer Academic Publishers, 1992.]。

(1) 発話様式  $H_{Style(N)}(\omega)$ 

加法性雑音  $N$  に依存する発話様式（しゃべり方、発話速度、発話の大きさ、Lombard効果など）に固有の周波数伝達特性である。Lombard効果とは、加法性雑音が存在する環境下で発声をする場合に、静寂な環境下とは異なって、無意識のうちに発声スペクトルが変形する現象のことをいう。文献[Y. Chen, "Cepstral Domain Talker Stress Compensation for Robust Speech Recognition," IEEE Trans. ASSP, Vol.36, No.4, pp.433-439, 1988.]では、ソフトなしゃべり方の場合は、1 kHz 以下のエネルギーが強く、1 kHz 以上のエネルギーが弱いという特性があること、一方で、大きな声、早口、叫び声、Lombard効果の場合は、逆の特性を持つことが示されている。

(2) 空間伝達特性  $H_{Trans}(\omega)$ 

口からマイクまでの空間的な周波数伝達特性を表す。

(3) マイク特性  $H_{Mic}(\omega)$ 

マイクなどの入力系の電氣的な周波数伝達特性を表す。

一般に、音声と雑音の線形スペクトル領域での加法性が成り立つとすると、時刻  $t$  における周波数  $\omega$  での観測スペクトル  $O(\omega; t)$  は、

$$O(\omega; t) = H_{Mic}(\omega) \cdot [H_{Trans}(\omega) \cdot \{H_{Style(N)}(\omega) \cdot (H_{Person}(\omega) \cdot S^*(\omega; t))\} + N(\omega; t) + E(\omega; t)] \quad (13)$$

でモデル化できる[J. H. L. Hansen, B. D. Womack, and L. M. Arslan, "A Source Generator Based Production Model for Environmental Robustness in Speech Recognition," Proc. ICSLP 94, Yokohama, Japan, pp.1003-1006, 1994.]. ここで、 $N(\omega; t)$  は音源未知の加法性雑音スペクトルを表し、 $E(\omega; t)$  は音源既知の加法性雑音スペクトルを表す。

4 種類の乗法性歪みの内、 $H_{Mic}(\omega)$  はあらかじめ測定可能であるが、 $H_{Person}(\omega)$ 、 $H_{Style(N)}(\omega)$ 、 $H_{Trans}(\omega)$  を、実環境において音声認識システムのユ

ーザーに負荷をかけることなく分離して測定することは困難であると考えられる。また、例えば加法性雑音  $N(\omega; t)$ 、 $E(\omega; t)$  が存在しないとしても、観測スペクトルの長時間平均として(10)式と同様に求められた時不変フィルターのゲインには、上記4種類の乗法性歪みの混在が避けられない。そこで、改めて乗法性歪み  $H^*(\omega)$ 、加法性雑音  $\tilde{N}(\omega; t)$ 、 $\tilde{E}(\omega; t)$  を、それぞれ

$$H^*(\omega) = H_{Mic}(\omega) \cdot H_{Trans}(\omega) \cdot H_{Style(N)}(\omega) \cdot H_{Person}(\omega) \quad (14)$$

$$\tilde{N}(\omega; t) = H_{Mic}(\omega) \cdot N(\omega; t) \quad (15)$$

$$\tilde{E}(\omega; t) = H_{Mic}(\omega) \cdot E(\omega; t) \quad (16)$$

と定義すると、(13)式を以下のように簡単化できる。

$$O(\omega; t) = H^*(\omega) \cdot S^*(\omega; t) + \tilde{N}(\omega; t) + \tilde{E}(\omega; t) \quad (17)$$

一方、(17)式を変形すると、

$$S^*(\omega; t) = \frac{\tilde{O}(\omega; t) - \tilde{N}(\omega; t) - \tilde{E}(\omega; t)}{H^*(\omega)} \quad (18)$$

が得られる。不特定話者音素モデルを観測されたスペクトルではなく、

(11)式により正規化されたスペクトルを用いて作成しておけば、観測スペ

クトル  $O(\omega; t)$  に対し、実環境における  $\tilde{N}(\omega; t)$ 、 $\tilde{E}(\omega; t)$ 、 $H^*(\omega)$  の除去を行っ

て、正規化音声スペクトル  $S^*(\omega; t)$  の推定値を求めることにより、頑健な音

声認識システムを実現できると考えられる。 $\tilde{N}(\omega; t)$ 、 $\tilde{E}(\omega; t)$  の除去に関して

は、NLMS-VAD法とCSS法を組み合わせる方法が有効であることを先述した。

音声認識システムにおいては、音響パラメータとして、通常、スペクトルの代わりに、ケプストラムが用いられる。ケプストラムは、スペクトルの対数値に逆離散コサイン変換(DCT: Discrete Cosine Transform)を施したものであるとして定義される。ケプストラムはスペクトルに比べて、少ないパラメータ数で同等の音声認識性能が得られるためよく用いられる。

正規化音声スペクトル  $S^*(\omega; t)$  に対する乗法性歪み  $H^*(\omega)$  の除去に関しては、次のE-CMN(Exact Cepstrum Mean Normalization)法が有効であることが既に示されている[M. Shozakai, S. Nakamura and K. Shikano, "A Non-Iterative Model-Adaptive E-CMN/PMC Approach for Speech Recognition in Car Environments," Proc. Eurospeech, Rhodes, Greece, pp.287-290, 1997.]. E-CMN法は次の2つのステップから構成される。推定ステップ: 話者毎に音声/非音声フレームで別々にケプストラム平均を求める。フレーム  $t$  における次数  $i$  のケプストラムを  $c(i, t)$  と表すとき、フレーム  $t$  における音声フレームのケプストラム平均  $\bar{C}_{speech}(i, t)$  は、例えば(19)式により求めることができる。

$$\bar{C}_{speech}(i, t) = \begin{cases} \eta \cdot \bar{C}_{speech}(i, t-1) + (1-\eta) \cdot c(i, t) & \text{if frame } t \text{ is speech} \\ \bar{C}_{speech}(i, t-1) & \text{otherwise} \end{cases} \quad (19)$$

ここで、 $\eta$  は音声フレームのケプストラム平均を求める際の平滑化係数であり、1.0よりも小さく1.0に近い値を設定すればよい。また、フレーム  $t$  における非音声フレームのケプストラム平均  $\bar{C}_{nonspeech}(i, t)$  は、例えば(20)式により求めることができる。ここで、 $\eta$  は非音声フレームのケプスト

ラム平均を求める際の平滑化係数であり、1.0よりも小さく1.0に近い値を設定すればよい。

$$\tilde{C}_{nonspeech}(i,t) = \begin{cases} \eta \cdot \tilde{C}_{nonspeech}(i,t-1) + (1-\eta) \cdot c(i,t) & \text{if frame } t \text{ is not speech} \\ C_{nonspeech}(i,t-1) & \text{otherwise} \end{cases}$$

(20)

音声フレームのケプストラム平均は、乗法性歪み  $H^*(\omega)$  のケプストラム表現であり、話者に依存する。一方、非音声フレームのケプストラム平均は、マイク特性などの入力系の乗法性歪み  $H_{Mic}(\omega)$  に依存する。

正規化ステップ：観測ケプストラム  $c(i,t)$  から、(21)式に従って、音声フレームにおいては音声フレームのケプストラム平均  $\tilde{C}_{speech}(i,t)$  を、非音声フレームにおいては非音声フレームのケプストラム平均  $\tilde{C}_{nonspeech}(i,t)$  を話者毎に引き、正規化ケプストラム  $\hat{C}(i,t)$  を求めることにより、観測スペクトルを正規化する。

$$\hat{C}(i,t) = \begin{cases} c(i,t) - \tilde{C}_{speech}(i,t) & \text{if frame } t \text{ is speech} \\ c(i,t) - \tilde{C}_{nonspeech}(i,t) & \text{otherwise} \end{cases} \quad (21)$$

尚、E-CMN法に関する発明について本願出願人は、既に日本国出願を済ませている（特願平09-051578号、1997年3月11日）。

E-CMN法は、様々な乗法性歪みの積を10単語程度の少量の音声から、音声区間のケプストラム平均として推定し、それを入力ケプストラムから引くという方法である。不特定話者音素モデルを観測されたスペクト

ルから求められたケプストラムではなく、E-CMN法により正規化されたケプストラムを用いて作成しておくことにより、様々な乗法性歪みを一括して補正することが可能であることが明らかにされている。

最後に、音源既知および音源未知の加法性雑音、乗法性歪みが存在する実環境におけるロバストな音声認識手法として、NLMS-VAD法、CSS法、E-CMN法を組み合わせる手法について述べる。図12に本組み合わせ手法に従って構成した演算回路のブロック図を示す。まず、第1の回路101においてNLMS-VAD法により、入力音声から音源既知の加法性雑音 $\hat{E}(\omega; t)$ が除去された波形信号が生成される。次に、第2の回路102においてこの波形信号に対して、フーリエ変換が施された後、CSS法により音源未知の加法性雑音 $\hat{N}(\omega; t)$ が除去されたスペクトルの時系列が生成される。

さらに、第3の回路103においてこのスペクトルの時系列はケプストラムの時系列に変換され、E-CMN法により正規化されたケプストラムの時系列に変換される。最後に、第4の回路104においてケプストラムの時系列は、公知のビタビアルゴリズムにより、不特定話者用に予め作成された音素モデルと照合され、音声認識結果が出力される。

この際に使用される不特定話者用音素モデルは、先述のようにE-CMN法により正規化されたケプストラムを用いて作成されている必要がある。尚、E-CMN法で必要な音声フレーム・非音声フレームの区別は、NLMS-VAD法に組み込まれたVADの結果をそのまま用いればよい。

次に、本組み合わせ手法の効果についてまとめる。単一指向性マイクを2000ccの自動車の運転席サンバイザーに設置し、男性2名女性2名が各々好みの位置にセットした運転席に座って発声した520単語（AT

R音声データベースCセット)の音声(データ1)を収録した。音声区間の前後に250msずつの無音区間が付属するように手動で切り出しを行った。また、アイドリング、時速60km、時速100kmの走行状態で、5種類の音楽ソース(ポップス、ジャズ、ロック、クラシック、落語)を順番にカーオーディオで再生し、音楽ソースの左右チャンネルを混合した信号(データ2)と、マイク入力信号(データ3)のペアを同時に録音した。カーオーディオの出力ボリュームは、各走行状態毎に男性1名が快適と感じる音量にセットした。

アイドリング、時速60km、時速100kmでの音響エコーのマイクへの最大入力レベルはそれぞれ、60.7dBA、65.9dBA、70.6dBAであった。データ1とデータ3を計算機上で加算して評価データを作成した。データ2は、NLMS-VAD法の遠端入力として使用した。認識には、環境独立な54音素の不特定話者用Tied-MixtureHMMモデル(40名の音声データから作成)を用いた。分析条件は8kHzサンプリング、フレーム長32ms、フレームシフト10msで、特徴パラメータは、10次MFCC、10次 $\Delta$ MFCC、 $\Delta$ エネルギーであり、HMMモデルが共有する正規分布の数は、それぞれ256、256、64である。

不特定話者、520単語の認識タスクで、アイドリング、時速60km、時速100kmの走行状態で、スピーカー出力音が存在しない場合(w/o Speaker Out)、スピーカー出力音が存在するが、NLMS-VAD法を行わない場合(w/ Speaker Out w/o NLMS-VAD)、スピーカー出力音が存在し、NLMS-VAD法を行う場合(w/ Speaker Out w/ NLMS-VAD)の認識性能(5種類の音楽ソースの平均)およびRREを図13に示す。

いずれの走行状態でも80%以上のRREが得られた。また、NLMS-VAD法でも回復できない誤認識率は、アイドリング、時速60km、

時速 100 km でそれぞれ 0.7%、2.1%、1.8% と僅かであり、先の組み合わせ法の有効性が確認できた。

(第 1 実施形態)

第 1 実施形態の音声処理装置の回路構成を図 14 に示す。以下に述べる個々の手段は周知の回路、たとえば、デジタル回路、コンピュータやデジタルプロセッサの演算処理により実現する回路を使用できるので、当業者であれば、図 13 により音声処理装置を製作できるであろう。まず、サンプル毎の処理について説明する。音源信号 11 はスピーカー出力信号 13 として、スピーカーから出力される。供給手段 a a 7 は、現時点のインパルス応答 (FIR フィルターの係数) を保持し、生成手段 a a 4 に現時点のインパルス応答 16 を供給する。

音源信号 11 は生成手段 a a 4 に送られ、生成手段 a a 4 で FIR フィルターにより、疑似音響エコー信号 15 が生成される。除去手段 a a 5 において、疑似音響エコー信号 15 はマイク入力信号 12 から減じられ、音源エコーキャンセル信号 14 が生成される。更新手段 a a 3 において、音源信号 11 と音源エコーキャンセル信号 14 と供給手段 a a 7 が保持する現時点のインパルス応答 16 から、インパルス応答を更新し、更新インパルス応答 17 を生成して、それを供給手段 a a 7 に供給する。

供給手段 a a 7 は後述の音声判定フラグ 18 が OFF の間にのみ、更新手段 a a 3 から供給された更新インパルス応答 17 を新しい現時点のインパルス応答として保持し、音声判定フラグ 18 が ON の間は、更新手段 a a 3 から供給された更新インパルス応答 17 を棄却する。除去手段 a a 5 で生成された音響エコーキャンセル信号 14 は、判定手段 a a 1 にも送られる。

次に、フレーム毎の処理について説明する。サンプル毎に判定手段 a a



1 に送られた音響エコーキャンセル信号 1 4 は判定手段 a a 1 においてバッファに格納され、1 フレーム分たまった段階で判定手段 a a 1 でマイク入力手段（不図示）に音声が存在するかどうかの判定が行われ、音声判定フラグ 1 8 が出力される。肯定判定の場合（音声の存在が検出される場合）は、音声判定フラグの値は O N であるという。否定判定の場合（音声の存在が検出されない場合）は、音声判定フラグの値は O F F であるという。

この判定処理では、音響エコーキャンセル信号 1 4 の時間情報および周波数情報を利用して、音声信号がマイク入力信号 1 2 に含まれているかどうかを判定する。判定の方法としては、音源未知の加法性雑音に重畳した音声を検出できるアルゴリズム VAD (Voice Activity Detection) を用いると良い。例えば、Recommendation GSM 06.32 では、信号の L P C 分析後の残差エネルギーをしきい値と比較することにより音声の検出を行うが、音源未知の加法性雑音のエネルギーレベルに合わせてしきい値を変動させることができるため、音源未知の加法性雑音と音声を分離することができる。

しきい値の適応化を行う際には、スペクトルの定常性（周波数情報）とピッチ性（時間情報）を利用している。上記判定処理の結果、音声判定フラグ 1 8 が O F F の場合は、制御手段 a a 6 が供給手段 a a 7 に保持された現時点インパルス応答 1 6 を取り出し、所望のインパルス応答として保存手段 a a 2 に格納する。

一方、音声判定フラグ 1 8 が O N の場合は、供給手段 a a 7 が保持しているインパルス応答が所望の値から乖離している可能性があるため、制御手段 a a 6 は保存手段 a a 2 から保存インパルス応答を 1 つ取り出し、供給手段 a a 7 が保持しているインパルス応答に上書きする。保存手段 a a 2 は、1 つ以上のインパルス応答を保存できる F I F O であればよい。

### (第2実施形態)

図15に第2実施形態の基本構成を示す。まず、第1の実施形態で説明した図14の構成を持つ音声処理装置100は、音源信号11およびマイク入力信号12を用いて、マイク入力信号12に含まれる音響エコーをキャンセルし、音響エコーキャンセル信号14を生成する。次に、スペクトル計算手段bb1において、一定フレーム周期毎にフーリエ変換により音響エコーキャンセル信号14のスペクトル21を算出する。

スペクトル21はケプストラム計算手段bb4に送られ、一定フレーム毎のケプストラム24に変換される。照合手段bb5では、一定フレーム毎のケプストラム24を用いて、照合が行われ、認識結果25が出力される。照合に当たっては、周知の隠れマルコフモデルによる手法または周知の動的計画法による手法または周知のニューラルネットワークによる手法のいずれかを用いればよい。

### (第3実施形態)

図16に第3実施形態の基本構成を示す。図15の第2実施形態と同様の個所には同一の符号を付している。まず、図14の構成を持つ音声処理装置100は、音源信号11およびマイク入力信号12を用いて、マイク入力信号12に含まれる音響エコーをキャンセルし、音響エコーキャンセル信号14を生成する。次に、スペクトル計算手段bb1において、一定フレーム周期毎にフーリエ変換により音響エコーキャンセル信号14のスペクトル21を算出する。スペクトル21は、スペクトル平均計算手段bb2に送られ、一定フレーム毎に(8)式によりスペクトル平均22が求められる。

一方、スペクトル計算手段bb1で求められたスペクトル21は、スペクトル平均減算手段bb3に供給され、(9)式によりスペクトル平均22が

減じられ、雑音除去スペクトル 2 3 が求められる。雑音除去スペクトル 2 3 はケプストラム計算手段 b b 4 に送られ、一定フレーム毎のケプストラム 2 4 に変換される。照合手段 b b 5 では、一定フレーム毎のケプストラム 2 4 を用いて、照合が行われ、認識結果 2 5 が出力される。照合に当たっては、周知の隠れマルコフモデルによる手法または周知の動的計画法による手法または周知のニューラルネットワークによる手法のいずれかを用いればよい。

#### (第 4 実施形態)

図 1 7 に第 4 実施形態の基本構成を示す。図 1 7 において第 2 または第 3 実施形態と同様の個所には同一の符号を付している。まず、図 1 4 の構成を持つ音声処理装置 1 0 0 は、音源信号 1 1 およびマイク入力信号 1 2 を用いて、マイク入力信号 1 2 に含まれる音響エコーをキャンセルし、音響エコーキャンセル信号 1 4 を生成する。次に、スペクトル計算手段 b b 1 において、一定フレーム周期毎にフーリエ変換により音響エコーキャンセル信号 1 4 のスペクトル 2 1 を算出する。スペクトル 2 1 は、スペクトル平均計算手段 b b 2 に送られ、一定フレーム毎に(8)式によりスペクトル平均 2 2 が求められる。

一方、スペクトル計算手段 b b 1 で求められたスペクトル 2 1 は、スペクトル平均減算手段 b b 3 に供給され、(9)式によりスペクトル平均 2 2 が減じられ、雑音除去スペクトル 2 3 が求められる。雑音除去スペクトル 2 3 はケプストラム計算手段 b b 4 に送られ、一定フレーム毎のケプストラム 2 4 に変換される。

ケプストラム 2 4 はケプストラム平均計算手段 c c 1 に送られ、そこでケプストラム平均 3 1 が求められる。ケプストラム平均の算出に当たっては、例えば(19)式、(20)式を用いればよい。次に、ケプストラム平均減算手

段cc2において、ケプストラム24からケプストラム平均31が減算し、正規化ケプストラム32が算出される。尚、減算に当たっては、(21)式を用いればよい。照合手段bb5では、一定フレーム毎の正規化ケプストラム32を用いて、照合が行われ、認識結果25が出力される。照合に当たっては、周知の隠れマルコフモデルによる手法または周知の動的計画法による手法または周知のニューラルネットワークによる手法のいずれかを用いればよい。

#### (第5実施形態)

図18に第5実施形態の基本構成を示す。図18において第2、第3または第4実施形態と同様の個所には同一の符号を付している。まず、図14の構成を持つ音声処理装置100は、音源信号11およびマイク入力信号12を用いて、マイク入力信号12に含まれる音響エコーをキャンセルし、音響エコーキャンセル信号14を生成する。次に、スペクトル計算手段bb1において、一定フレーム周期毎にフーリエ変換により音響エコーキャンセル信号14のスペクトル21を算出する。スペクトル21は、ケプストラム計算手段bb4に送られ、一定フレーム毎のケプストラム24に変換される。ケプストラム24はケプストラム平均計算手段cc1に送られ、そこでケプストラム平均31が求められる。ケプストラム平均の算出に当たっては、例えば(19)式、(20)式を用いればよい。

次に、ケプストラム平均減算手段cc2において、ケプストラム24からケプストラム平均31を減算し、正規化ケプストラム32が算出される。尚、減算に当たっては、(21)式を用いればよい。照合手段bb5では、一定フレーム毎の正規化ケプストラム32を用いて、照合が行われ、認識結果25が出力される。照合に当たっては、周知の隠れマルコフモデルによる手法または周知の動的計画法による手法または周知のニューラルネット

ワークによる手法のいずれかを用いればよい。

(第6実施形態)

図19に第6実施形態の基本構成を示す。図16においても第2、第3第4または第5実施形態と同様の個所には同一の符号を付している。まず、図14の構成を持つ音声処理装置100により、音源信号11およびマイク入力信号12を用いて、マイク入力信号12に含まれる音響エコーをキャンセルし、音響エコーキャンセル信号14を生成する。次に、スペクトル計算手段bb1において、一定フレーム周期毎にフーリエ変換により音響エコーキャンセル信号14のスペクトル21を算出する。スペクトル21は、ケプストラム計算手段bb4に送られ、一定フレーム毎のケプストラム24に変換される。

ケプストラム24はケプストラム平均計算手段cc1に送られ、そこでケプストラム平均31が求められる。ケプストラム平均の算出に当たっては、例えば(19)式、(20)式を用いればよい。次に、ケプストラム平均減算手段cc2において、ケプストラム24からケプストラム平均31を減算し、正規化ケプストラム32が算出される。尚、減算に当たっては、(21)式を用いればよい。照合手段bb5では、一定フレーム毎の正規化ケプストラム32を用いて、照合が行われ、認識結果25が出力される。照合に当たっては、周知の隠れマルコフモデルによる手法または周知の動的計画法による手法または公知のニューラルネットワークによる手法のいずれかを用いればよい。

## 請求の範囲

1. 音響エコーの伝達経路を模擬する現時点のインパルス応答および音源信号に基づき疑似音響エコー信号を生成する生成手段と、

現時点のインパルス応答を保持し、前記生成手段に供給する供給手段と、

マイク入力信号から前記疑似音響エコー信号を減算することにより音響エコー成分を除去し、音響エコーキャンセル信号を生成する除去手段と、

前記音源信号と前記音響エコーキャンセル信号と前記供給手段が保持している現時点のインパルス応答を用いて継続的にインパルス応答を更新し、更新されたインパルス応答を前記供給手段に供給する更新手段と、

前記音響エコーキャンセル信号の時間情報および周波数情報を利用して、マイク入力信号に音声が含まれているか否かをフレーム毎に判定する判定手段と、

1つ以上のインパルス応答を保存する保存手段と、

前記判定手段の判定結果が否定判定のフレームでは前記供給手段が保持している現時点のインパルス応答を前記保存手段に保存し、肯定判定のフレームでは前記保存手段に保存されているインパルス応答の1つを取り出して、前記供給手段に供給する制御手段と

を具えたことを特徴とする音声処理装置。

2. 請求の範囲第1項に記載の音声処理装置において、前記音響エコーがキャンセルされた後の信号を音声認識に用いることを特徴とする音声処理装置。

3. 請求の範囲第2項に記載の音声処理装置において、前記音響エコーがキャンセルされた後の信号から、フーリエ変換により各フレーム毎にスペ

クトルを求める手段と、当該得られたスペクトルに基づき各フレーム毎に連続的にスペクトル平均を求める手段と、当該得られたスペクトル平均を前記音響エコーがキャンセルされた後の信号から各フレーム毎に計算されたスペクトルから連続的に減算することにより、音源未知の加法性雑音を除去する手段とをさらに具えたことを特徴とする音声処理装置。

4. 請求の範囲第2項に記載の音声処理装置において、前記音響エコーがキャンセルされた後の信号から、フーリエ変換により各フレーム毎にスペクトルを求める手段と、当該得られたスペクトルから各フレーム毎に連続的にスペクトル平均を求める手段と、当該得られたスペクトル平均を前記音響エコーがキャンセルされた後の信号から各フレーム毎に計算されたスペクトルを連続的に減算する手段と、当該減算する手段により、音源未知の加法性雑音が除去されたスペクトルからケプストラムを求める手段と、当該得られたケプストラムの音声フレームのケプストラム平均および非音声フレームのケプストラム平均を話者毎に別々に求める手段と、話者毎に音声フレームのケプストラムからはその話者の音声フレームのケプストラム平均を減算し、非音声フレームのケプストラムからはその話者の非音声フレームのケプストラム平均を減算して、マイク特性や口からマイクまでの空間伝達特性に依存する乗法性歪みを補正する手段をさらに具えたことを特徴とする音声処理装置。

5. 請求の範囲第2項に記載の音声処理装置において、前記音響エコーがキャンセルされた後の信号から、フーリエ変換により各フレーム毎にスペクトルを求める手段と、当該得られたスペクトルからケプストラムを求める手段と、当該得られたケプストラムの音声フレームのケプストラム平均および非音声フレームのケプストラム平均を話者毎に別々に求める手段と、話者毎に音声フレームのケプストラムからはその話者の音声フレームのケ

プストラム平均を減算し、非音声フレームのケプストラムからはその話者の非音声フレームのケプストラム平均を減算することにより、マイク特性や口からマイクまでの空間伝達特性に依存する乗法性歪みを補正する手段とをさらに具えたことを特徴とする。

6. フーリエ変換により各フレーム毎にスペクトルを求める手段と、

当該得られたスペクトルからケプストラムを求める手段と、

当該得られたケプストラムの音声フレームのケプストラム平均および非音声フレームのケプストラム平均を話者毎に別々に求める手段と、

話者毎に音声フレームのケプストラムからはその話者の音声フレームのケプストラム平均を減算し、非音声フレームのケプストラムからはその話者の非音声フレームのケプストラム平均を減算することによりマイク特性や口からマイクまでの空間伝達特性に依存する乗法性歪みを補正する手段と

を具えたことを特徴とする音声処理装置。

7. 音響エコーの伝達経路を模擬する現時点のインパルス応答および音源信号に基づき疑似音響エコー信号を生成する生成ステップと、

現時点のインパルス応答を保持し、前記生成手段に供給する供給ステップと、

マイク入力信号から前記疑似音響エコー信号を減算することにより音響エコー成分を除去し、音響エコーキャンセル信号を生成する除去ステップと、

前記音源信号と前記音響エコーキャンセル信号と前記供給ステップで保持している現時点のインパルス応答を用いて継続的にインパルス応答を更新し、更新されたインパルス応答を前記供給ステップに供給する更新ステップと、



前記音響エコーキャンセル信号の時間情報および周波数情報を利用して、マイク入力信号に音声が含まれているか否かをフレーム毎に判定する判定ステップと、

1つ以上のインパルス応答を保存する保存ステップと、

前記判定ステップの判定結果が否定判定のフレームでは前記供給ステップが保持している現時点のインパルス応答を前記保存ステップで保存し、肯定判定のフレームでは前記保存ステップで保存されているインパルス応答の1つを取り出して、前記供給ステップに供給する制御ステップと

を具えたことを特徴とする音声処理方法。

8. 請求の範囲第7項に記載の音声処理方法において、前記音響エコーがキャンセルされた後の信号を音声認識に用いることを特徴とする音声処理方法。

9. 請求の範囲第8項に記載の音声処理方法において、前記音響エコーがキャンセルされた後の信号から、フーリエ変換により各フレーム毎にスペクトルを求めるステップと、当該得られたスペクトルに基づき各フレーム毎に連続的にスペクトル平均を求めるステップと、当該得られたスペクトル平均を前記音響エコーがキャンセルされた後の信号から各フレーム毎に計算されたスペクトルを連続的に減算することにより、音源未知の加法性雑音を除去するステップとをさらに具えたことを特徴とする音声処理方法。

10. 請求の範囲第8項に記載の音声処理方法において、前記音響エコーがキャンセルされた後の信号から、フーリエ変換により各フレーム毎にスペクトルを求めるステップと、当該得られたスペクトルから各フレーム毎に連続的にスペクトル平均を求めるステップと、当該得られたスペクトル平均を前記音響エコーがキャンセルされた後の信号から各フレーム毎に計算されたスペクトルから連続的に減算することにより、音源未知の加法性

雑音を除去するステップと、当該加法性雑音が除去されたスペクトルからケプストラムを求めるステップと、当該得られたケプストラムの音声フレームのケプストラム平均および非音声フレームのケプストラム平均を話者毎に別々に求めるステップと、話者毎に音声フレームのケプストラムからはその話者の音声フレームのケプストラム平均を減算し、非音声フレームのケプストラムからはその話者の非音声フレームのケプストラム平均を減算して、マイク特性や口からマイクまでの空間伝達特性に依存する乗法性歪みを補正するステップをさらに具えたことを特徴とする音声処理方法。

11. 請求の範囲第8項に記載の音声処理方法において、前記音響エコーがキャンセルされた後の信号から、フーリエ変換により各フレーム毎にスペクトルを求めるステップと、当該得られたスペクトルからケプストラムを求めるステップと、当該得られたケプストラムの音声フレームのケプストラム平均および非音声フレームのケプストラム平均を話者毎に別々に求めるステップと、話者毎に音声フレームのケプストラムからはその話者の音声フレームのケプストラム平均を減算し、その話者の非音声フレームのケプストラムからは非音声フレームのケプストラム平均を減算することにより、マイク特性や口からマイクまでの空間伝達特性に依存する乗法性歪みを補正するステップとをさらに具えたことを特徴とする音声処理方法。

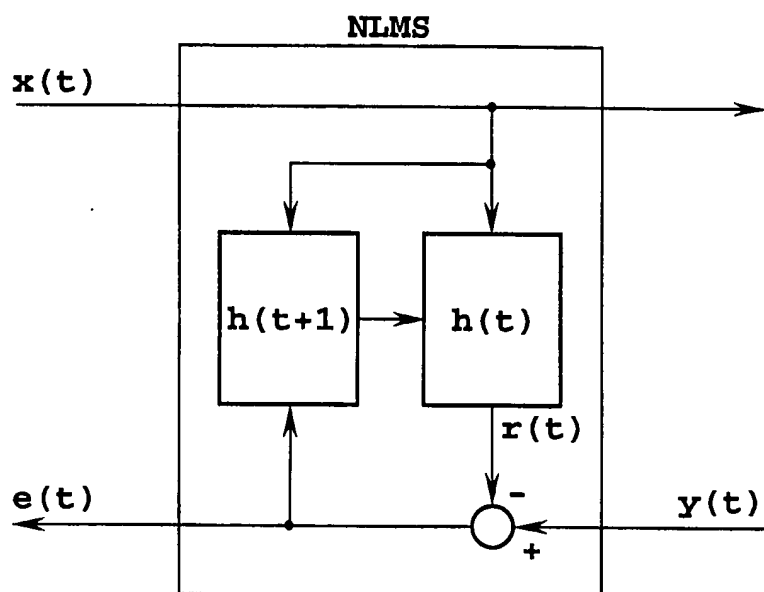
12. フーリエ変換により各フレーム毎にスペクトルを求めるステップと、  
当該得られたスペクトルからケプストラムを求めるステップと、  
当該得られたケプストラムの音声フレームのケプストラム平均および非音声フレームのケプストラム平均を話者毎に別々に求めるステップと、

話者毎に音声フレームのケプストラムからはその話者の音声フレームのケプストラム平均を減算し、非音声フレームのケプストラムからはその話者の非音声フレームのケプストラム平均を減算することによりマイク特性

や口からマイクまでの空間伝達特性に依存する乗法性歪みを補正するステップと

を具えたことを特徴とする音声処理方法。

1/20

**FIG.1**

2/20

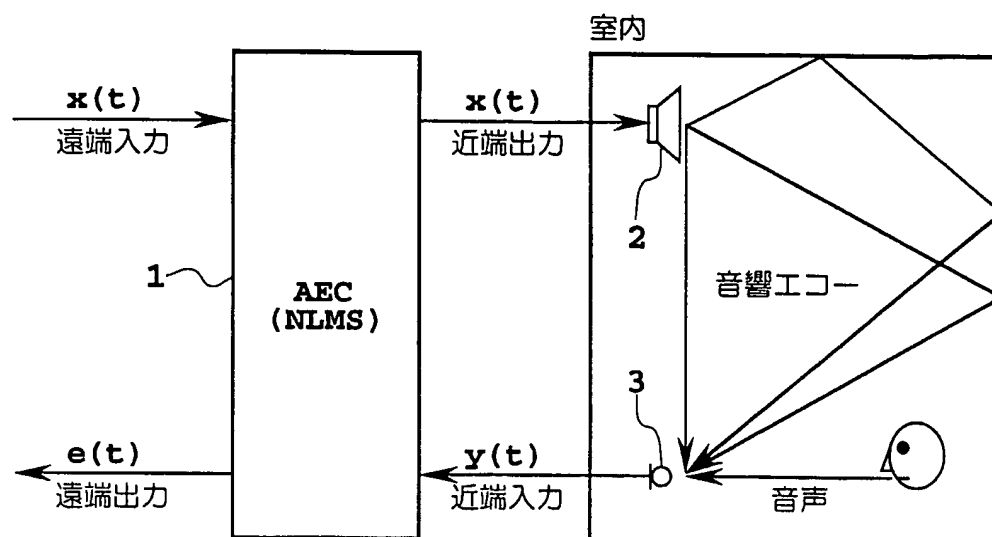


FIG.2

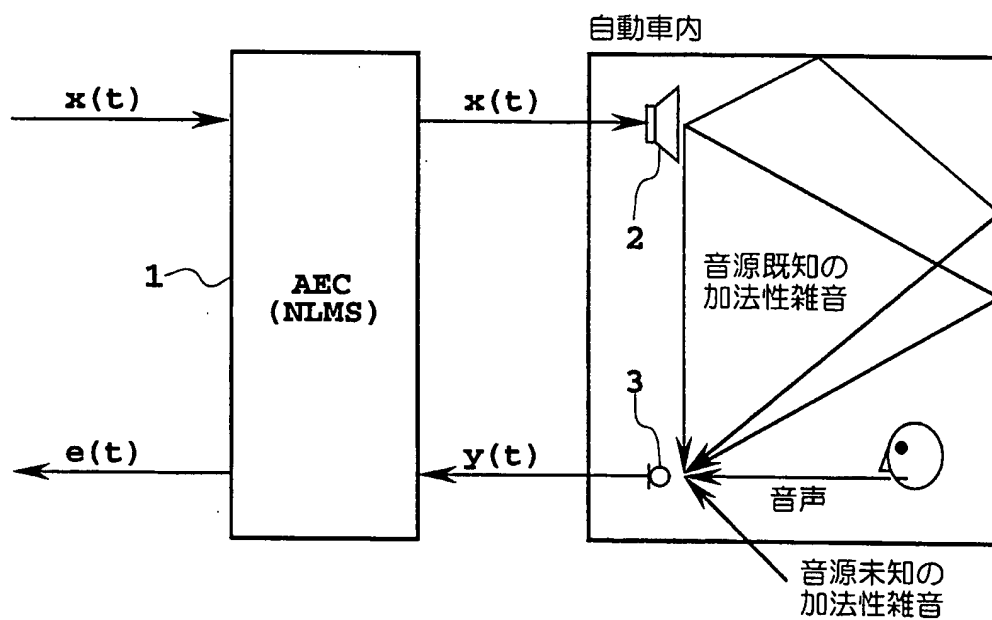
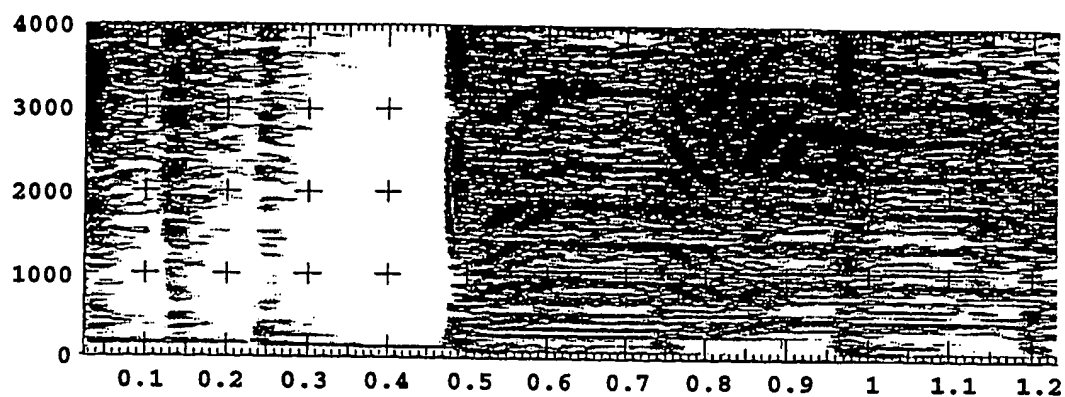
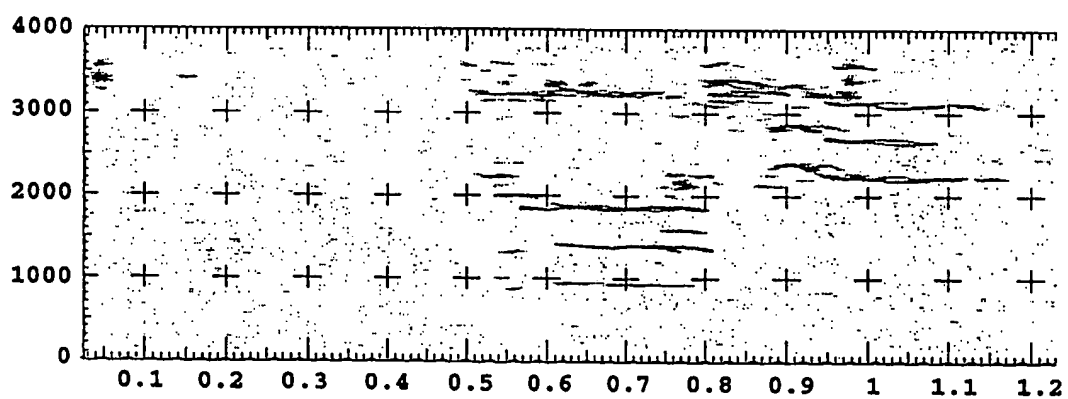
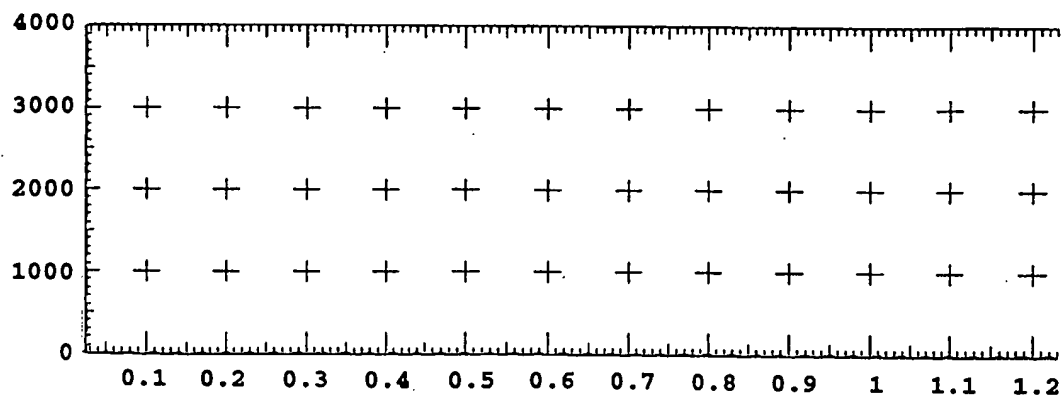
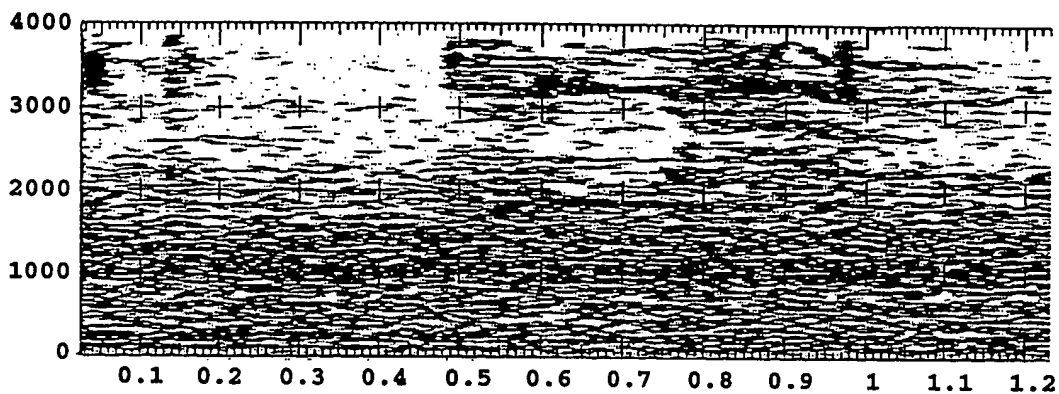
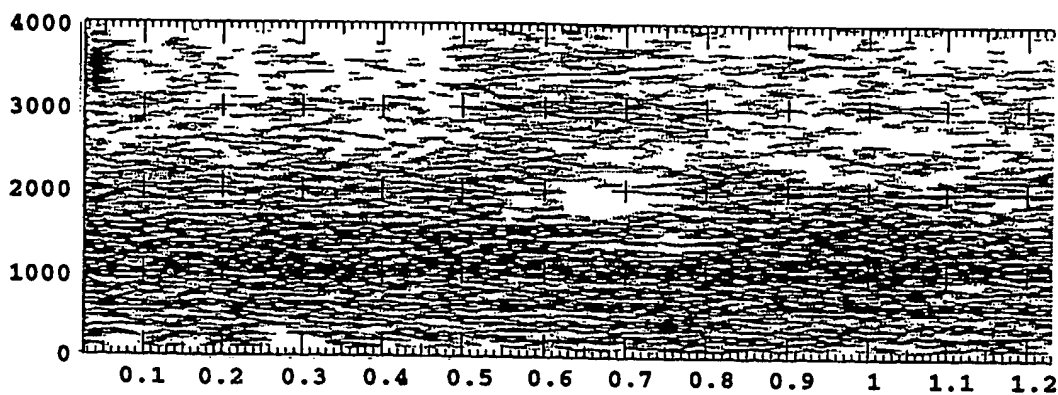


FIG.3

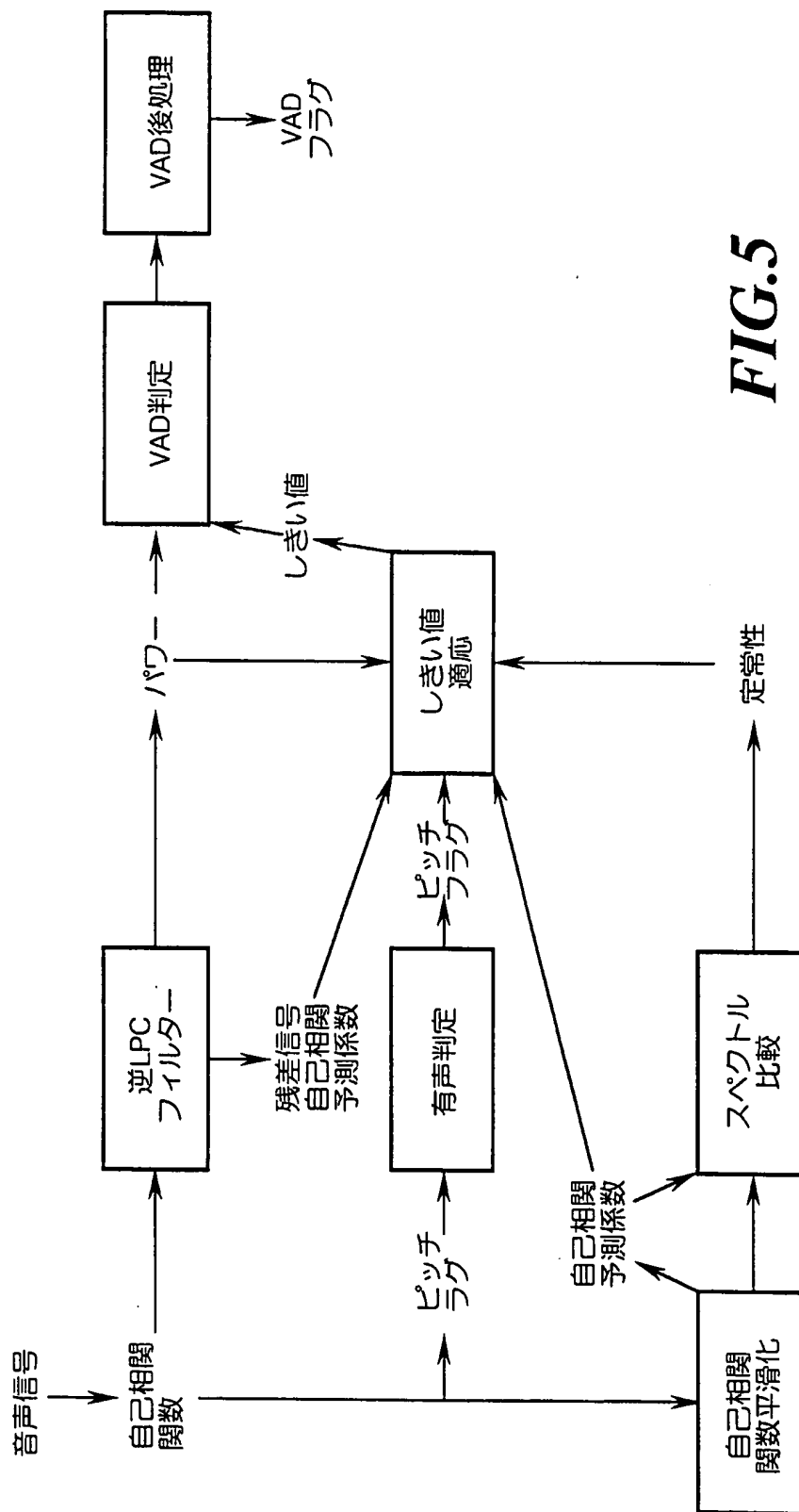
3/20

**FIG. 4A****FIG. 4B**

4/20

*FIG.4C**FIG.4D**FIG.4E*

**5 / 2 0**





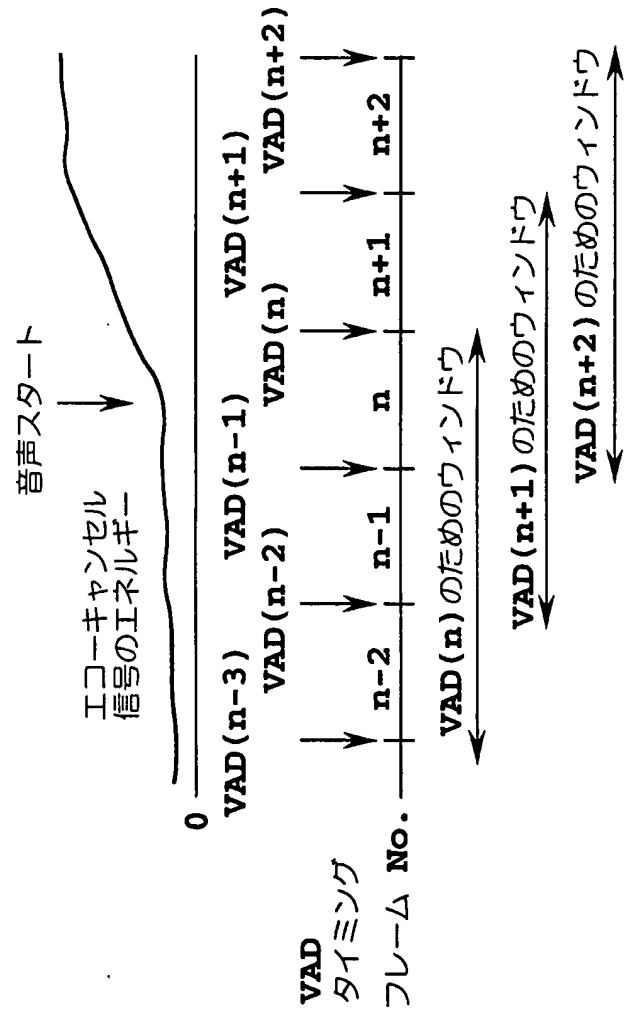
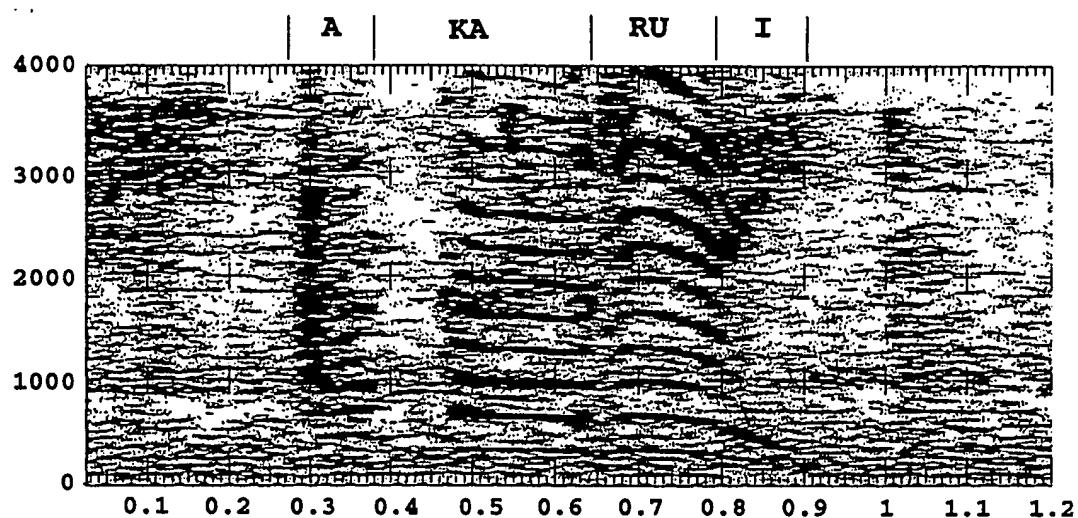
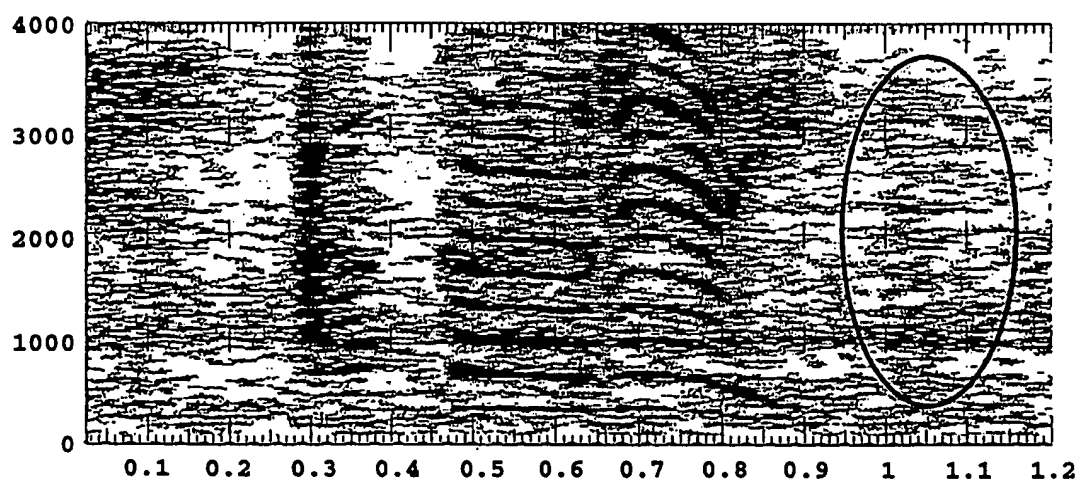
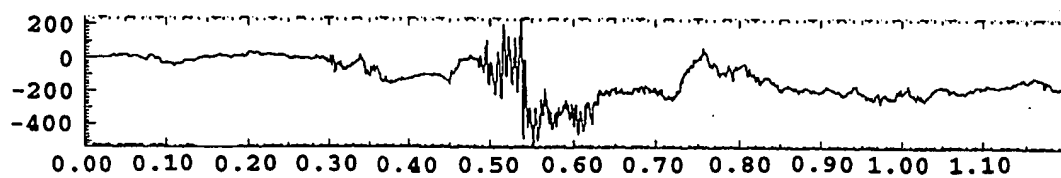
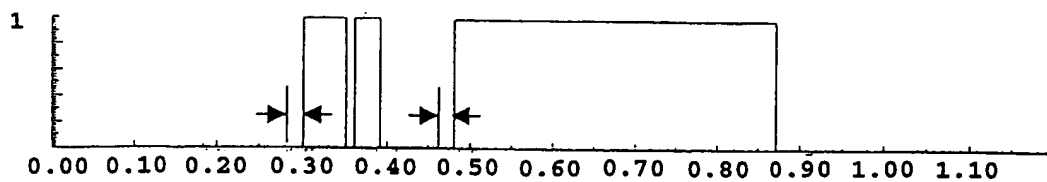
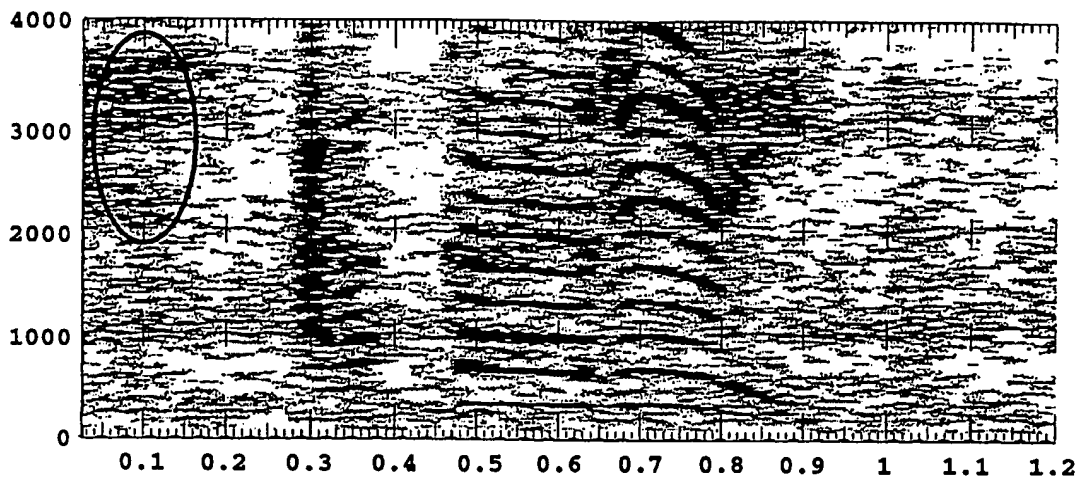
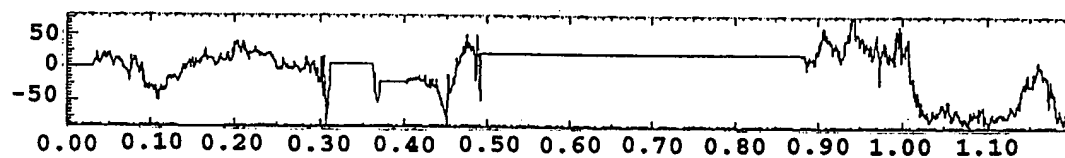
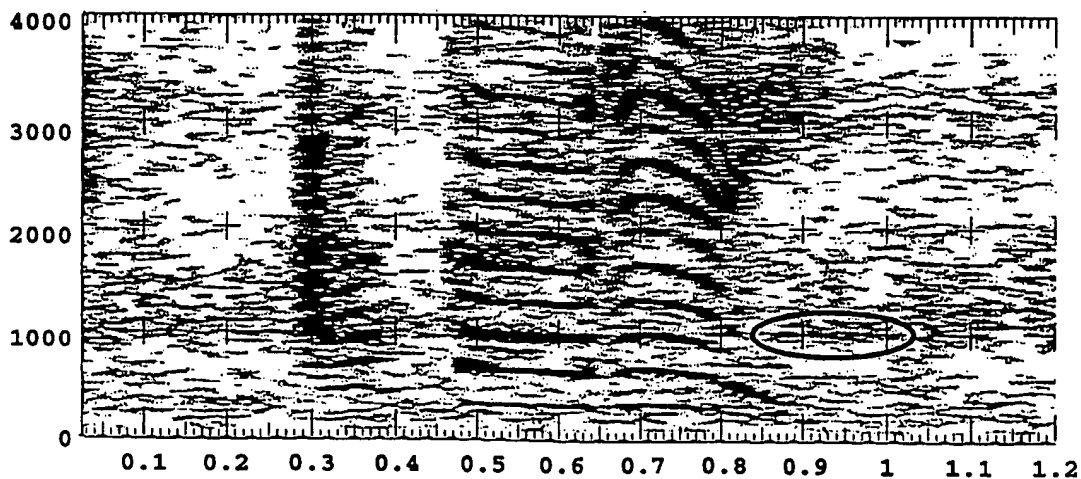


FIG.6

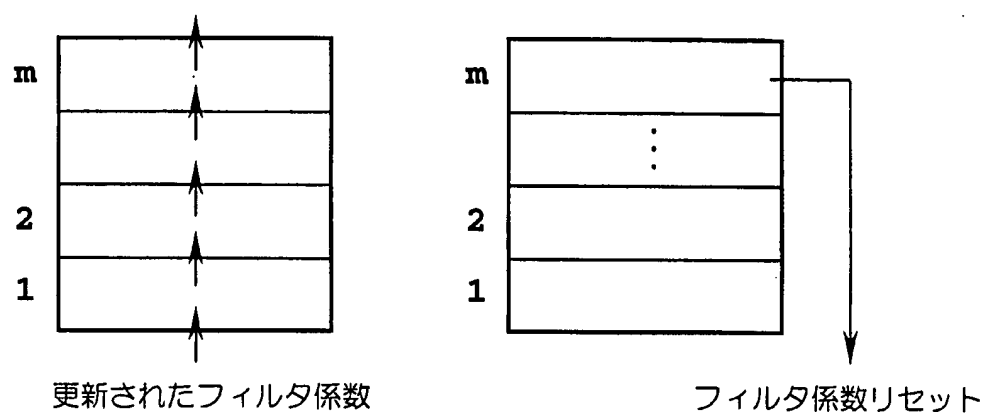
7/20

**FIG. 7A****FIG. 7B****FIG. 7C**

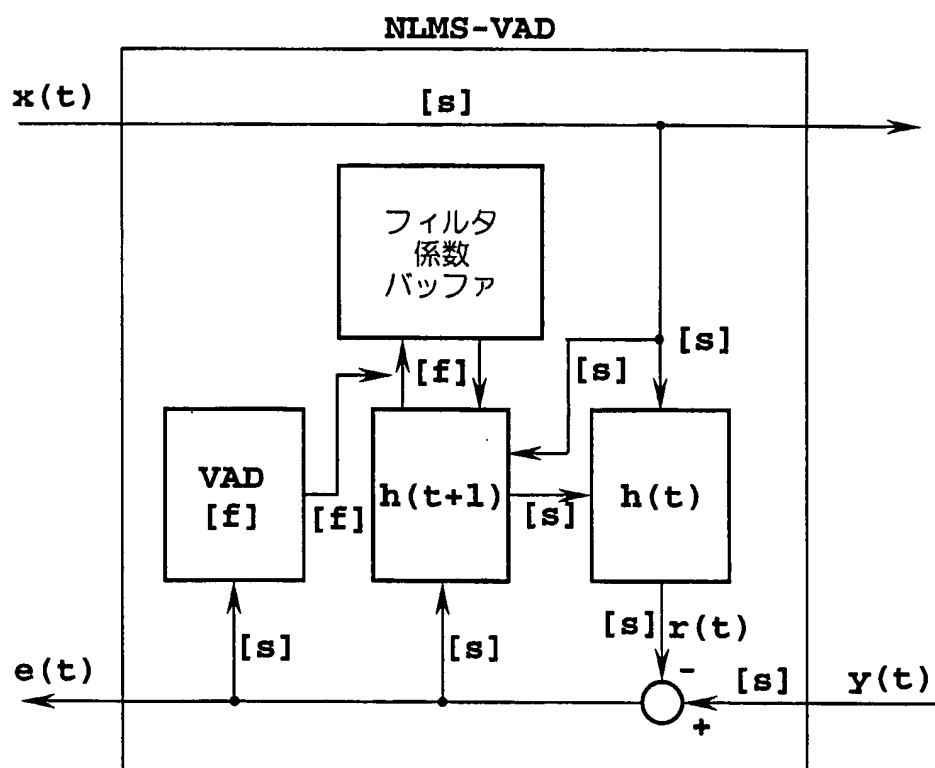
8/20

**FIG. 7D****FIG. 7E****FIG. 7F****FIG. 7G**

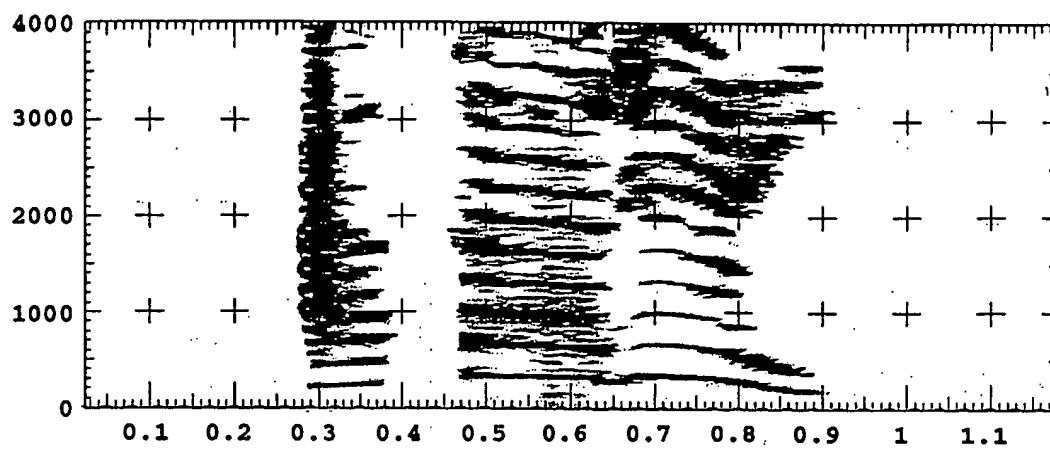
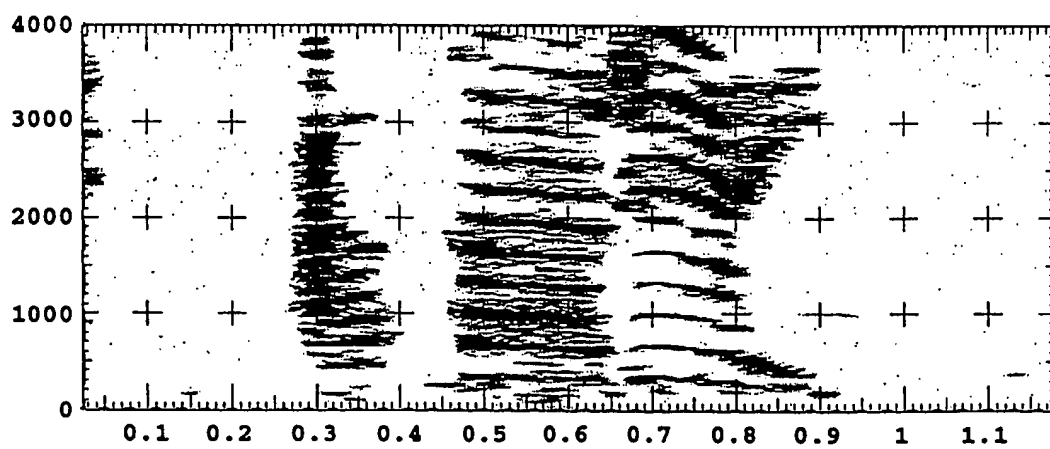
9/20

**FIG.8**

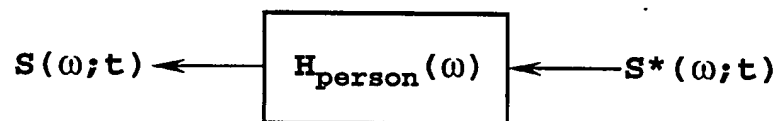
**10/20**



**FIG.9**

*11/20**FIG. 10A**FIG. 10B*

12/20

**FIG.11**

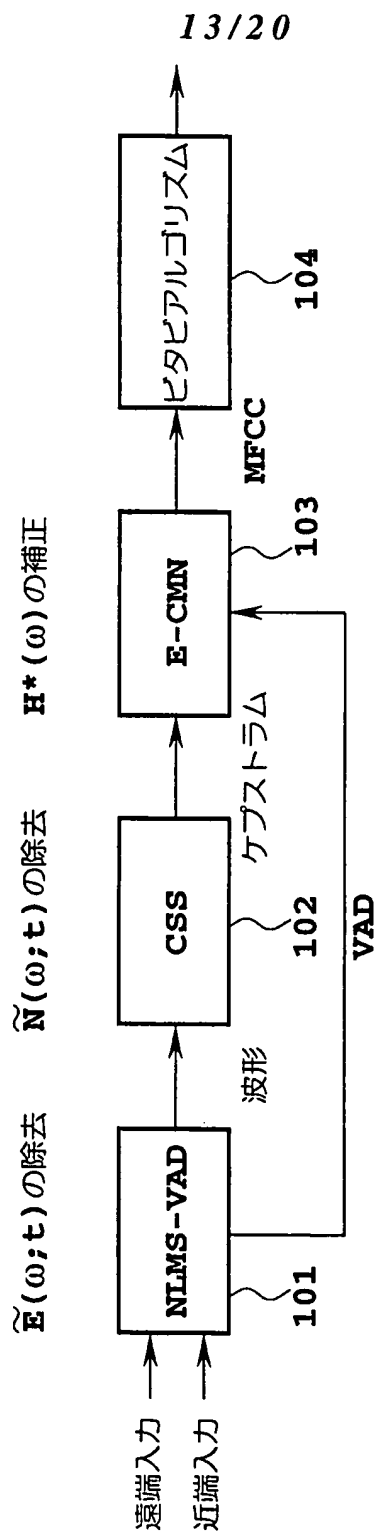


FIG.12



14/20

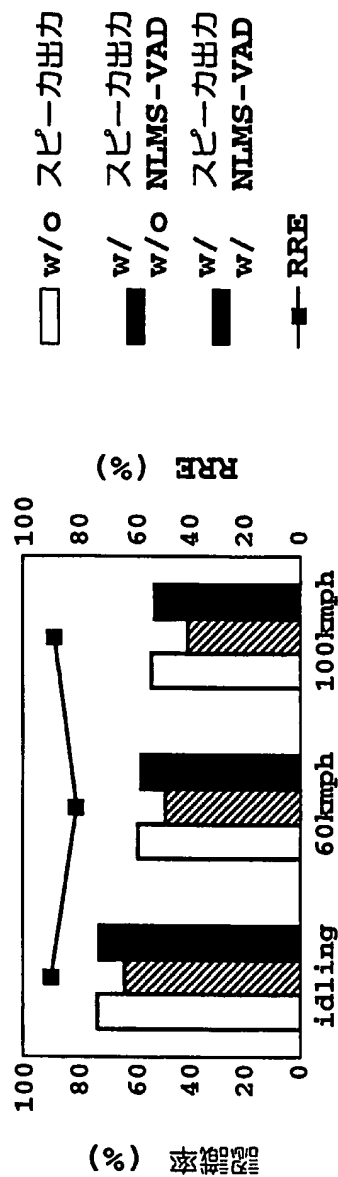


FIG.13

15/20

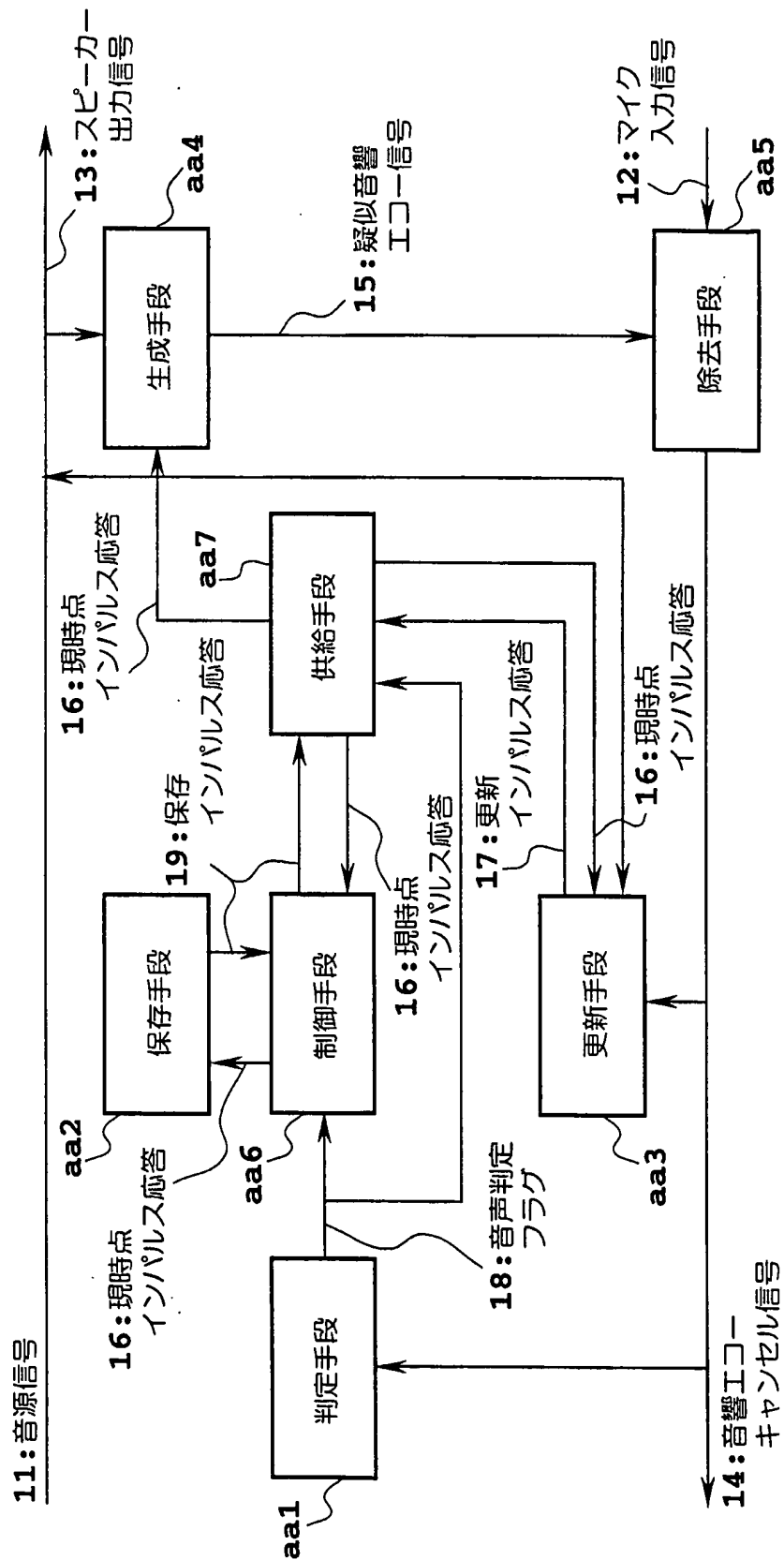


FIG.14

16/20

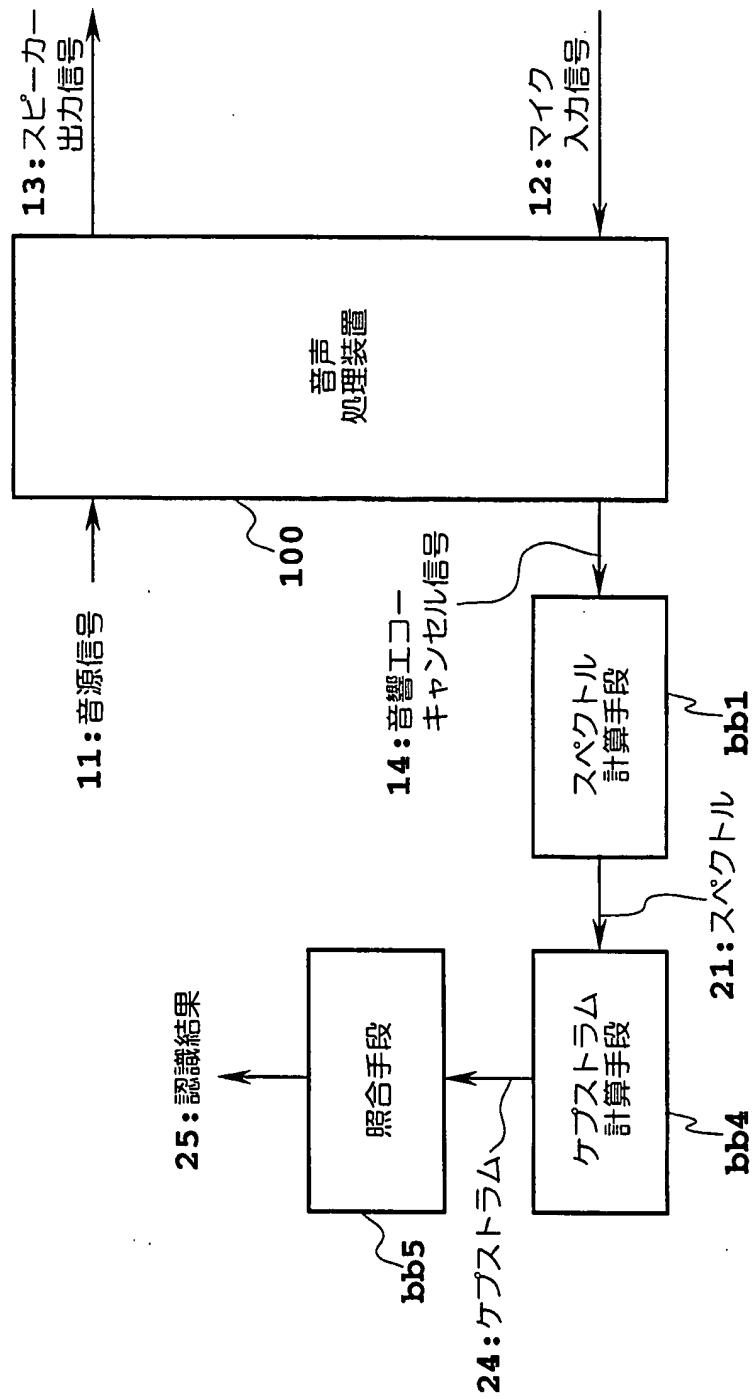


FIG.15

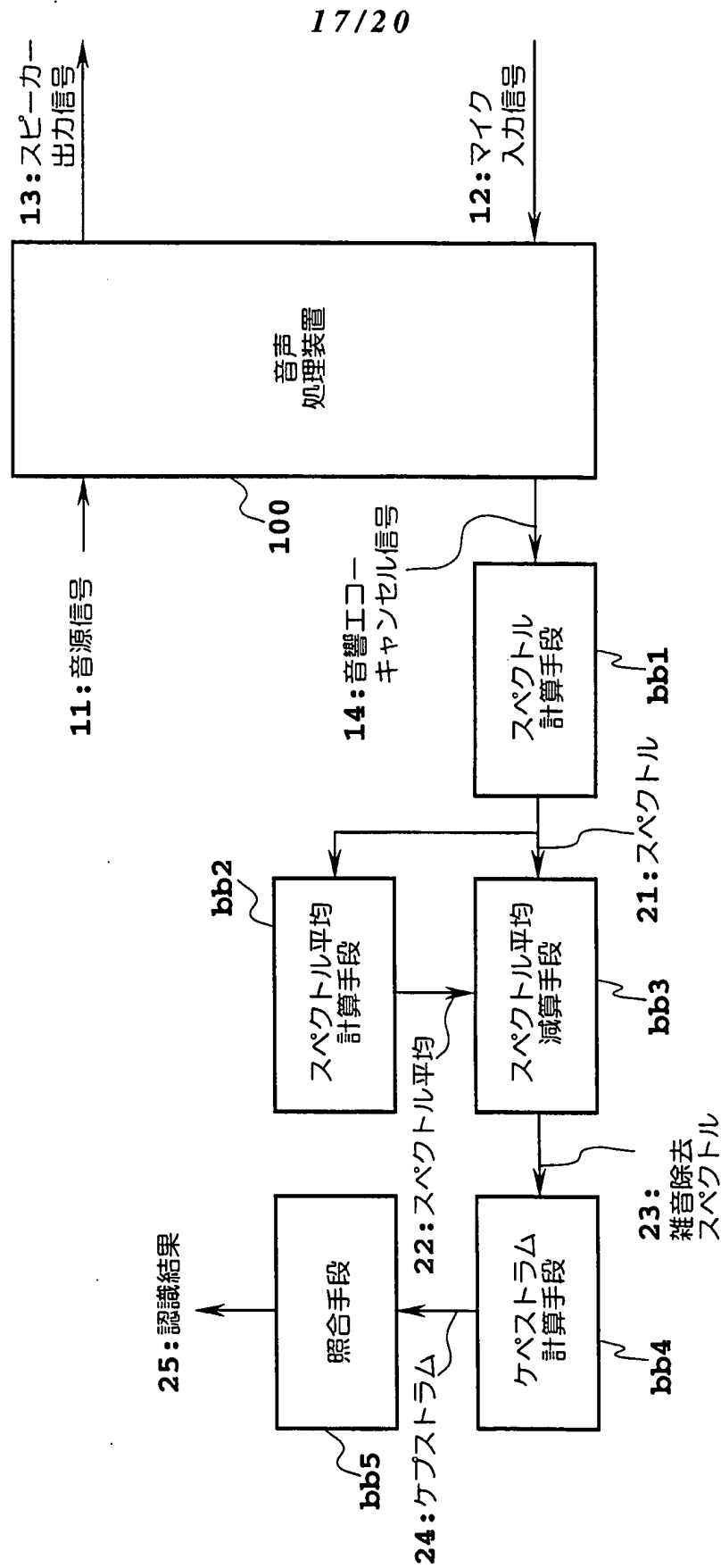


FIG.16

18/20

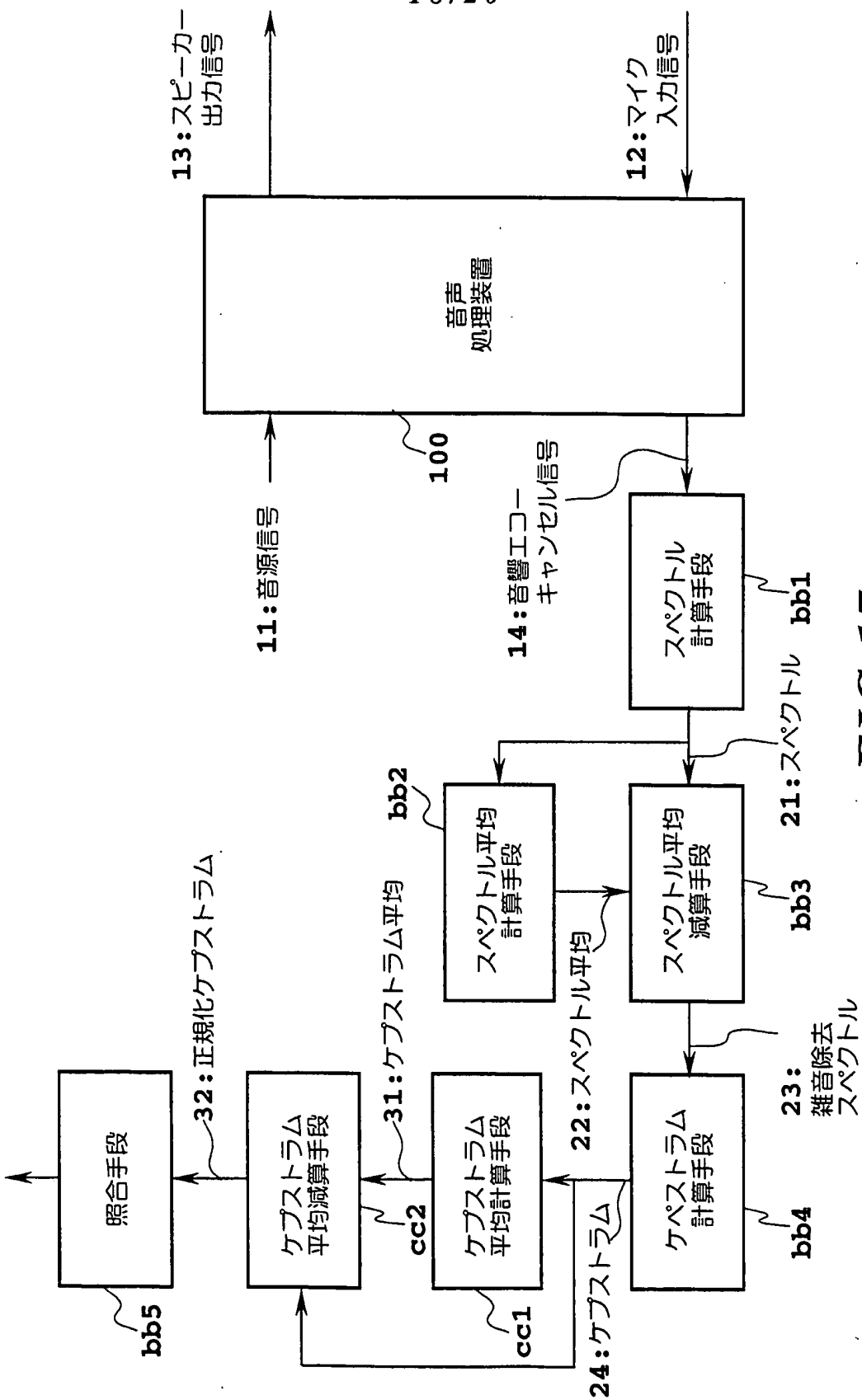


FIG.17

19/20

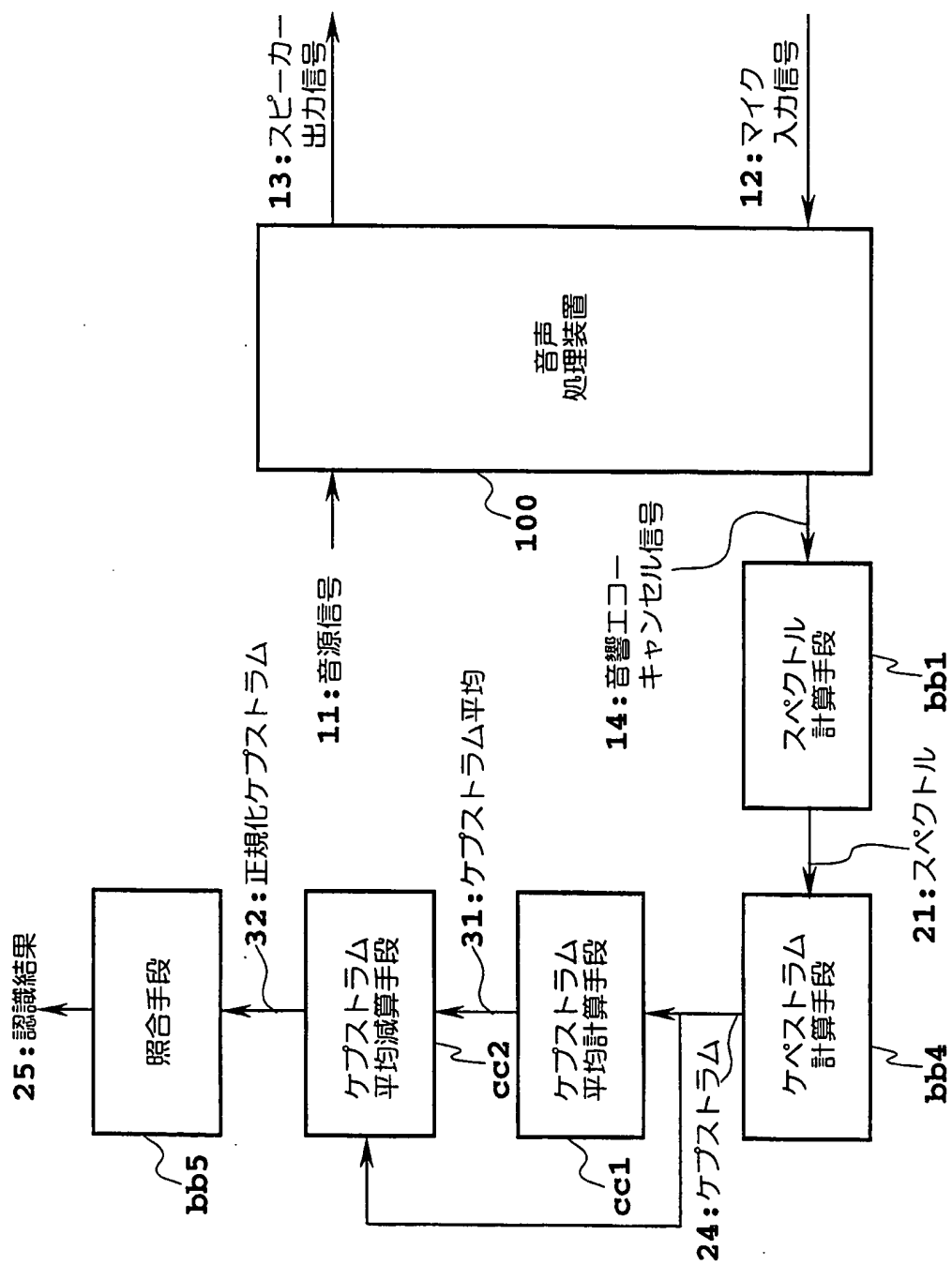


FIG. 18

20/20

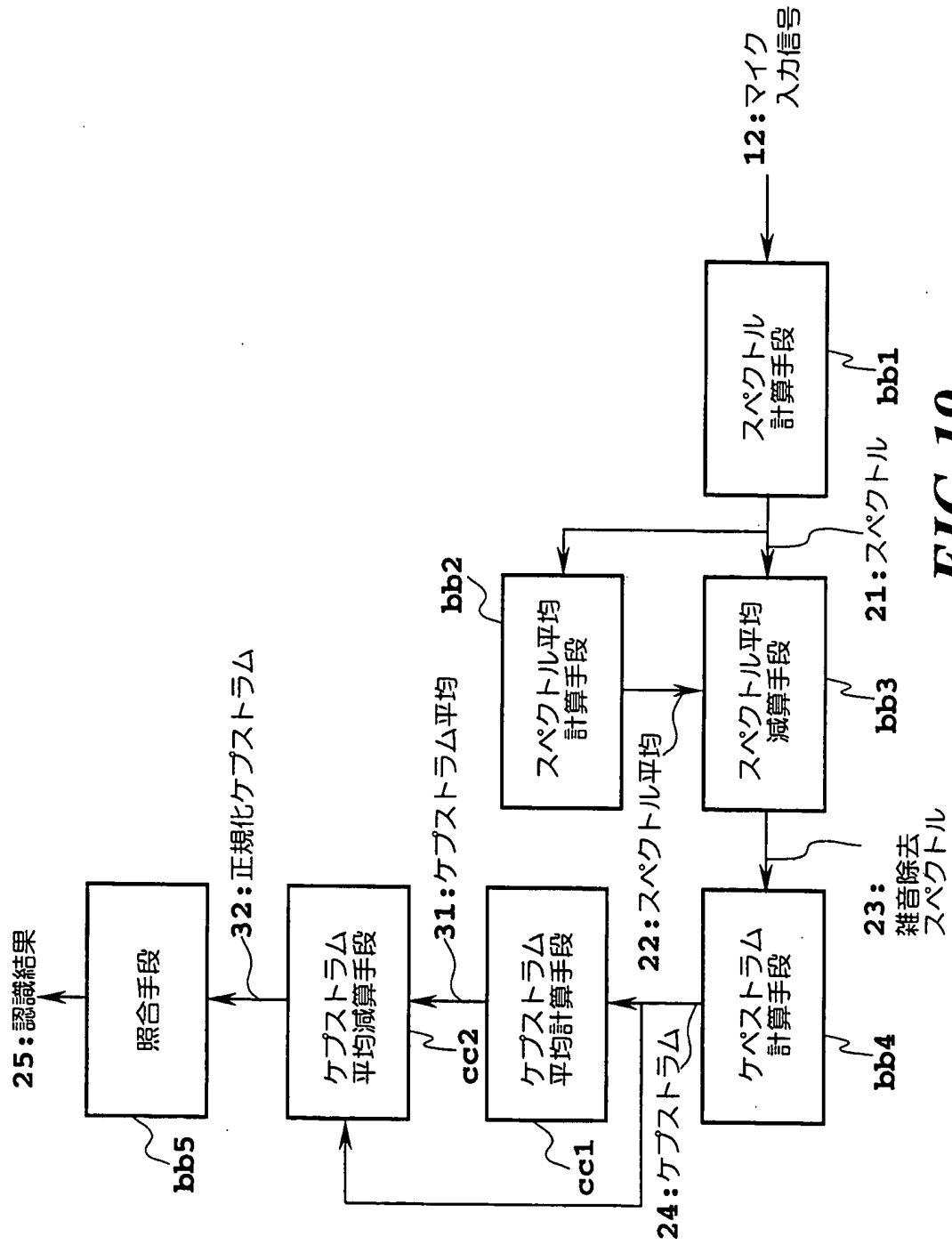


FIG. 19

## INTERNATIONAL SEARCH REPORT

International application No.

PCT/JP98/00915

A. CLASSIFICATION OF SUBJECT MATTER  
Int.Cl<sup>6</sup> H04R3/02

According to International Patent Classification (IPC) or to both national classification and IPC

## B. FIELDS SEARCHED

Minimum documentation searched (classification system followed by classification symbols)

Int.Cl<sup>6</sup> H04R3/02, H03H17/00, H04M1/16, H04B3/23

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched

Jitsuyo Shinan Koho 1926-1995

Kokai Jitsuyo Shinan Koho 1971-1995

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)

## C. DOCUMENTS CONSIDERED TO BE RELEVANT

Category*	Citation of document, with indication, where appropriate, of the relevant passages	Relevant to claim No.
X	JP, A, 63-18797 (Matsushita Electric Industrial Co., Ltd.), January 26, 1988 (26. 01. 88) (Family: none)	1, 7
Y	JP, A, 7-66757 (Nippon Telegraph and Telephone Corp.), March 10, 1995 (10. 03. 95) (Family: none)	2, 3, 8, 9

☐ Further documents are listed in the continuation of Box C.☐ See patent family annex.

\* Special categories of cited documents:

"A" document defining the general state of the art which is not considered to be of particular relevance

"E" earlier document but published on or after the international filing date  
"L" document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified)

"O" document referring to an oral disclosure, use, exhibition or other means

"P" document published prior to the international filing date but later than the priority date claimed

"T" later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention

"X" document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone

"Y" document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art

"&amp;" document member of the same patent family

Date of the actual completion of the international search  
June 2, 1998 (02. 06. 98)Date of mailing of the international search report  
June 16, 1998 (16. 06. 98)Name and mailing address of the ISA/  
Japanese Patent Office

Authorized officer

Facsimile No.

Telephone No.



## 国際調査報告

国際出願番号 PCT/J P 98/00915

## A. 発明の属する分野の分類 (国際特許分類 (IPC))

Int. cl.<sup>8</sup> H04R 3/02

## B. 調査を行った分野

## 調査を行った最小限資料 (国際特許分類 (IPC))

Int. cl.<sup>8</sup> H04R 3/02 H03H 17/00 H04M 1/16 H04B 3/23

## 最小限資料以外の資料で調査を行った分野に含まれるもの

日本国実用新案公法 1926-1995

日本国公開実用新案公法 1971-1995

## 国際調査で使用した電子データベース (データベースの名称、調査に使用した用語)

## C. 関連すると認められる文献

引用文献の カテゴリー*	引用文献名 及び一部の箇所が関連するときは、その関連する箇所の表示	関連する 請求の範囲の番号
X	J P, A, 63-18797, (松下電器産業株式会社), 26. 1月. 1988 (26. 01. 88), ファミリなし	1, 7
Y	J P, A, 7-66757, (日本電信電話株式会社), 10. 3月. 1995 (10. 03. 95), ファミリなし	2, 3, 8, 9

☐ C欄の続きにも文献が列挙されている。☐ パテントファミリーに関する別紙を参照。

## \* 引用文献のカテゴリー

「A」特に関連のある文献ではなく、一般的技術水準を示すもの

「E」先行文献ではあるが、国際出願日以後に公表されたもの

「L」優先権主張に疑義を提起する文献又は他の文献の発行日若しくは他の特別な理由を確立するために引用する文献 (理由を付す)

「O」口頭による開示、使用、展示等に言及する文献

「P」国際出願日前で、かつ優先権の主張の基礎となる出願

の日の後に公表された文献

「T」国際出願日又は優先日後に公表された文献であって出願と矛盾するものではなく、発明の原理又は理論の理解のために引用するもの

「X」特に関連のある文献であって、当該文献のみで発明の新規性又は進歩性がないと考えられるもの

「Y」特に関連のある文献であって、当該文献と他の1以上の文献との、当業者にとって自明である組合せによって進歩性がないと考えられるもの

「&amp;」同一パテントファミリー文献

国際調査を完了した日

02. 06. 98

国際調査報告の発送日

16.06.98

国際調査機関の名称及びあて先

日本国特許庁 (ISA/J P)

郵便番号100-8915

東京都千代田区霞が関三丁目4番3号

特許庁審査官 (権限のある職員)

新 宮 佳 典

5 H

7 5 2 5

電話番号 03-3581-1101 内線 3530